

(19) 日本国特許庁(JP)

(12) 特 許 公 報(B2)

(11) 特許番号

特許第4087097号
(P4087097)

(45) 発行日 平成20年5月14日(2008.5.14)

(24) 登録日 平成20年2月29日(2008.2.29)

(51) Int.Cl.

F I

G06F 12/00 (2006.01)
G06F 17/30 (2006.01)

G06F 12/00 501B
G06F 12/00 513Z
G06F 17/30 110C
G06F 17/30 240Z

請求項の数 14 (全 42 頁)

(21) 出願番号 特願2001-345523 (P2001-345523)
(22) 出願日 平成13年11月12日(2001.11.12)
(65) 公開番号 特開2003-150414 (P2003-150414A)
(43) 公開日 平成15年5月23日(2003.5.23)
審査請求日 平成15年12月24日(2003.12.24)

(73) 特許権者 000005108
株式会社日立製作所
東京都千代田区丸の内一丁目6番6号
(74) 代理人 100079108
弁理士 稲葉 良幸
(74) 代理人 100093861
弁理士 大賀 眞司
(72) 発明者 茂木 和彦
神奈川県川崎市麻生区王禅寺1099番地
株式会社日立製作所 システム開発研究
所内
(72) 発明者 大枝 高
神奈川県川崎市麻生区王禅寺1099番地
株式会社日立製作所 システム開発研究
所内

最終頁に続く

(54) 【発明の名称】 データベース管理システム情報を考慮したデータ再配置方法およびデータ再配置を行う計算機システム

(57) 【特許請求の範囲】

【請求項1】

データベース管理システムが稼動する少なくとも1台以上の計算機と、
少なくとも1つ以上のデータを記憶する物理記憶手段を有し、前記データベース管理システムにより管理されるデータベースデータを格納した少なくとも1台以上の記憶装置と、

前記計算機と前記記憶装置との間に接続され、前記計算機と前記記憶装置との間でデータの転送を制御する記憶制御手段と、

前記記憶装置における前記データの配置を管理するデータ位置管理サーバとを有する計算機システムにおけるデータの再配置方法において、

前記データ位置管理サーバが、前記計算機から前記データベース管理システムにより管理されるデータベースに関する情報、前記記憶制御手段からデータの記憶位置に関する情報及び前記記憶装置から前記物理記憶手段の稼動情報である物理記憶稼動情報を取得するステップと、

前記データ位置管理サーバが、前記データベースに関する情報、前記データの記憶位置に関する情報及び前記物理記憶稼動情報からなる取得情報に基づいて、前記データベースデータの配置を決定するステップと、

前記データ位置管理サーバから前記決定されたデータ配置を実現するためのデータ移動を前記記憶制御手段に指示するステップと、

前記記憶制御手段が、前記指示に従って、前記記憶装置内に格納された前記データの配

置を変更するステップとを有し、

前記データベースデータの配置を決定するステップでは、前記データの配置を決定する際に、前記取得情報に基づいて、同時にアクセスされる可能性が高い前記データベースデータの組を検出し、前記組を異なる前記物理記憶手段に配置するように、かつ、前記物理記憶手段における記憶位置まで定めたデータの配置を決定することを特徴とするデータの再配置方法。

【請求項 2】

前記データ位置管理サーバの機能が、所定の前記データベースが稼動する計算機上で実現されることを特徴とする請求項 1 記載のデータの再配置方法。

【請求項 3】

前記データ位置管理サーバの機能が、所定の前記記憶制御手段上で実現されることを特徴とする請求項 1 記載のデータの再配置方法。

【請求項 4】

前記データ位置管理サーバの機能が、所定の前記記憶装置上で実現されることを特徴とする請求項 1 記載のデータの再配置方法。

【請求項 5】

前記記憶制御手段が、前記計算機上で実施されるプログラムにより実現されることを特徴とする請求項 1 記載のデータの再配置方法。

【請求項 6】

前記データ位置管理サーバから前記記憶制御手段に決定されたデータの配置を指示する前に、前記決定されたデータの配置を管理者に提示し、管理者にデータの配置変更を実施するか確認することを特徴とする請求項 1 記載のデータの再配置方法。

【請求項 7】

前記データベースに関する情報は、前記データベース管理システムのスキーマにより定義される表・索引・ログを含むデータ構造に関する情報と、前記データベースデータを前記スキーマにより定義されるデータ構造毎に分類した前記記憶装置における記録位置に関する情報の少なくとも 1 つを含むことを特徴とする請求項 1 に記載のデータの再配置方法。

【請求項 8】

データベース管理システムが稼動する少なくとも 1 台以上の計算機と、
少なくとも 1 つ以上のデータを記憶する物理記憶手段を有し、前記データベース管理システムにより管理されるデータベースデータを格納した少なくとも 1 台以上の記憶装置と、

前記計算機と前記記憶装置との間に接続され、前記計算機と前記記憶装置との間でデータの転送を制御する記憶制御手段と、

前記記憶装置における前記データの配置を管理するデータ位置管理サーバとを有する計算機システムにおいて、

前記データ位置管理サーバは、前記計算機から前記データベース管理システムにより管理されるデータベースに関する情報、前記記憶制御手段からデータの記憶位置に関する情報及び前記記憶装置から前記物理記憶手段の稼動情報である物理記憶稼動情報を取得する情報取得手段と、

前記データベースに関する情報、前記データの記憶位置に関する情報及び前記物理記憶稼動情報からなる取得情報に基づいて、前記データベースデータの配置を決定する配置決定手段と、

前記データ位置管理サーバから前記記憶制御手段により決定されたデータの配置を実現するためのデータ移動を指示するデータ配置指示手段を有し、

前記記憶制御手段は、前記データ位置管理サーバからの指示に従って、前記複数の記憶装置内に格納された前記データの配置を変更するデータ配置変更手段を有し、

前記配置決定手段は、前記データの配置を決定する際に、前記取得情報に基づいて、同時にアクセスされる可能性が高い前記データベースデータの組を検出し、前記組を異なる

10

20

30

40

50

前記物理記憶手段に配置するように、かつ、前記物理記憶手段における記憶位置まで定めたデータの配置を決定することを特徴とする計算機システム。

【請求項 9】

前記データ位置管理サーバと所定の前記データベースが稼動する計算機が同一の計算機であることを特徴とする請求項 8 記載の計算機システム。

【請求項 10】

前記データ位置管理サーバと前記記憶制御手段が同一の装置であることを特徴とする請求項 8 記載の計算機システム。

【請求項 11】

前記データ位置管理サーバと所定の前記記憶装置が同一の装置であることを特徴とする請求項 8 記載の計算機システム。

10

【請求項 12】

前記記憶制御手段が、前記計算機上で実施されるプログラムにより実現されることを特徴とする請求項 8 記載の計算機システム。

【請求項 13】

前記配置決定手段により決定したデータの配置を管理者に提示する手段と、前記管理者からデータの配置変更の可否を取得する手段を有することを特徴とする請求項 8 記載の計算機システム。

【請求項 14】

前記データベースに関する情報は、前記データベース管理システムのスキーマにより定義される表・索引・ログを含むデータ構造に関する情報と、前記データベースデータを前記スキーマにより定義されるデータ構造毎に分類した前記記憶装置における記録位置に関する情報の少なくとも 1 つを含むことを特徴とする請求項 8 に記載の計算機システム。

20

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、データベース管理システムに関する。

【0002】

【従来の技術】

現在、データベース（DB）を基盤とする多くのアプリケーションが存在し、DBに関する一連の処理・管理を行うソフトウェアであるデータベース管理システム（DBMS）は極めて重要なものとなっている。特に、DBMSの処理性能はDBを利用するアプリケーションの性能も決定するため、DBMSの処理性能の向上は極めて重要である。

30

【0003】

DBの特徴の1つは、多大なデータを扱うことである。そのため、DBMSの実行環境の多くにおいては、DBMSが実行される計算機に対して大容量の記憶装置あるいは複数の記憶装置を接続し、記憶装置上にDBのデータを記憶する。そのため、DBに関する処理を行う場合に、記憶装置に対してアクセスが発生し、記憶装置におけるデータアクセス性能がDBMSの性能を大きく左右する。そのため、DBMSが稼動するシステムにおいて、記憶装置の性能最適化が極めて重要であり、特にDBMSが管理するデータの物理記憶装置に対する配置の最適化が重要である。

40

文献“Oracle 8i パフォーマンスのための設計およびチューニング，リリース 8.1”（日本オラクル社，部品番号 J00921-01）の第 20 章（文献 1）は、RDBMS である Oracle 8i における I/O のチューニングを述べている。その中で、RDBMS の内部動作のチューニングと共に、データの配置のチューニングに関連するものとして、ログファイルは他のデータファイルから分離した物理記憶装置に記憶すること、ストライプ化されたディスクにデータを記憶することによる負荷分散が効果があること、表のデータとそれに対応する索引データは異なる物理記憶装置に記憶すると効果があること、RDBMS とは関係ないデータを異なる物理記憶装置に記憶することが述べられている。

50

【 0 0 0 4 】

米国特許 6 0 3 5 3 0 6 (文献 2) は、 D B M S - ファイルシステム - ボリュームマネージャ - 記憶装置間のマッピングを考慮した性能解析ツールに関する技術を開示している。この性能解析ツールは、各レイヤにおけるオブジェクトの稼動状況を画面に表示する。このときに上記のマッピングを考慮し、その各オブジェクトに対応する他レイヤのオブジェクトの稼動状況を示す画面を容易に表示する機能を提供する。また、ボリュームマネージャレイヤのオブジェクトに関して、負荷が高い記憶装置群に記憶されているオブジェクトのうち、2 番目に負荷が高いオブジェクトを、もっとも負荷が低い記憶装置群に移動するオブジェクト再配置案を作成する機能を有している。

【 0 0 0 5 】

特開平 9 - 2 7 4 5 4 4 号公報 (文献 3) は、計算機がアクセスするために利用する論理的記憶装置を、実際にデータを記憶する物理記憶装置に配置する記憶装置において、前記論理的記憶装置の物理記憶装置への配置を動的に変更することにより記憶装置のアクセス性能を向上する技術を開示している。アクセス頻度が高い物理記憶装置に記憶されているデータの一部を前記の配置動的変更機能を用いて他の物理記憶装置に移動することにより、特定の物理記憶装置のアクセス頻度が高くならないようにし、これにより記憶装置を全体としてみたときの性能を向上させる。また、配置動的変更機能による高性能化処理の自動実行方法についても開示している。

【 0 0 0 6 】

特開 2 0 0 1 - 6 7 1 8 7 号公報 (文献 4) は、計算機がアクセスするために利用する論理的記憶装置を、実際にデータを記憶する物理記憶装置に配置し、前記論理的記憶装置の物理記憶装置への配置を動的に変更する機能を有する記憶装置において、論理的記憶装置の物理記憶装置への配置の変更案を作成する際に、物理記憶装置を属性の異なるグループに分割し、それを考慮した配置変更案を作成し、その配置変更案に従って自動的に論理的記憶装置の配置を変更する技術を開示している。配置変更案作成時に、物理記憶装置を属性毎にグループ化し、論理的記憶装置の配置先として、それが有している特徴にあった属性を保持している物理記憶装置のグループに属する物理記憶装置を選択する配置変更案を作成することにより、良好な配置変更案を作成する。

【 0 0 0 7 】

【 発明が解決しようとする課題 】

従来の技術には以下のような問題が存在する。

【 0 0 0 8 】

文献 1 で述べられているものは、管理者がデータの配置を決定する際に考慮すべき項目である。現在、1 つの記憶装置内に多数の物理記憶装置を有し、多数の計算機により共有されるものが存在する。この記憶装置においては、多くの場合、ホストが認識する論理的記憶装置を実際にデータを記憶する物理記憶装置の適当な領域に割り当てることが行われる。このような記憶装置を利用する場合、人間がすべてを把握することは困難であり、このような記憶装置を含む計算機システム側に何かしらのサポート機能が存在しなければ文献 1 に述べられている問題点を把握することすら困難となる。また、問題点を把握することができたとしても、計算機システム側にデータの移動機能が存在しない場合には、記憶装置上のデータの再配置のためにデータのバックアップとリストアが必要となり、その処理に多大な労力を必要とする。

【 0 0 0 9 】

文献 2 で述べられている技術では、ボリュームマネージャレイヤにおけるオブジェクトの稼動状況によるデータ再配置案を作成する機能を実現しているが、記憶装置から更に高いアクセス性能を得ようとする場合には D B M S レイヤにおけるデータの特徴を考慮して配置を決定する必要があるがその点の解決方法については何も述べていない。

【 0 0 1 0 】

文献 3 , 文献 4 で述べられている技術においては、記憶装置を利用するアプリケーションが利用するデータに関する特徴としてアクセス頻度とシーケンシャルアクセス性程度しか

10

20

30

40

50

考慮していないため、アプリケーションから見た場合に必ずしも良好な配置が実現できるわけではない。例えば、DBMSでは表データとそれに対応する索引データを同時にアクセスすることが多いが、これらのデータを同一の物理記憶装置に配置する可能性がある。この場合、物理記憶装置においてアクセス競合が発生し、記憶装置のアクセス性能が低下する可能性がある。

【0011】

本発明の第一の目的は、DBMSが管理するデータの特性を考慮して記憶装置におけるデータ記憶位置を決定することにより、DBMSに対してより好ましいアクセス性能特性を持つ記憶装置を保持する計算機システムを実現し、DBMSの性能を向上させることである。特に、複数の記憶装置を利用するDBシステムにおいて、各記憶装置間へアクセス要求が適切に分散化されるようにし、DBMSの処理性能を向上させる。

10

【0012】

本発明の第二の目的は、DBMSが稼動している計算機システムにおいて、DBMSの特性を考慮した記憶装置の良好なアクセス性能特性達成を目的としたデータ記憶位置再配置処理を実現することにより、計算機システムの性能に関する管理コストを削減することである。

【0013】

【課題を解決するための手段】

DBMSに関する情報を1箇所に集約させ、そこでデータの再配置案を決定し、そのデータ再配置案に従うようにデータの移動指示を発行することにより、DBMSに対してより好ましい性能特性を持つデータ配置を実現する。

20

【0014】

計算機システム内には、複数のDBMSが稼動し、複数の記憶装置が利用されている可能性がある。そのため、DBMSや記憶装置の構成その他の情報を一箇所に集約し、そこで集約した全情報を考慮してデータの再配置案を作成する。

【0015】

DBMSが管理するデータの配置決定方法として、以下のものを採用する。データ更新時に必ず書き込みが実行される更新ログを、他のデータと異なる物理記憶装置に配置して相互干渉しないようにすることによりDBMSに対して良好な性能特性を得ることができる。同時にアクセスされる可能性が極めて高い表データとそれに対応する索引データを異なる物理記憶装置に配置することによりDBMSに対して良好な性能特性を得ることができる。DBMSに関する情報を利用して、データがシーケンシャルにアクセスされる場合のアクセス順序を予測し、その構造を保持するように物理記憶装置に記憶する。

30

【0016】

現在、計算機のオペレーティングシステム(OS)・データ転送経路中のスイッチ・記憶装置内部においてデータの記憶位置を変更する技術が存在する。データの記憶位置の変更はそれらの技術を用いて実現する。そこで、前記の項目を考慮して決定されたデータ再配置案に従って、データの記憶位置の変更を管理する部分に対してデータ配置変更の指示を発行する。

【0017】

【発明の実施の形態】

以下、本発明の実施の形態を説明する。なお、これにより本発明が限定されるものではない。

<第一の実施の形態>

本実施形態では、DBMSが実行される計算機と記憶装置がスイッチを用いて接続された計算機システムにおいて、データの記憶位置の管理を行う計算機が存在し、そこで計算機システム内のデータの記憶位置の管理を行う。本実施例におけるスイッチは、記憶装置から提供される記憶領域を組み合わせて仮想的な記憶装置を作成する機能を有する。また、記憶装置においても、記憶装置内部でデータの記憶位置を動的に変更する機能を有する。

【0018】

40

50

データ記憶位置管理を実施する計算機は、DBMSに関する情報、データの記憶位置のマッピングに関する情報、記憶装置の構成情報を取得し、それらを用いて好適なデータ再配置案を作成する。スイッチと記憶装置に対して作成したデータ配置を指示し、それらのデータ再配置機能を用いてそのデータ再配置案に従ったデータ配置を実現する。

【0019】

図1は、本発明の第一の実施の形態における計算機システムの構成図である。本実施の形態における計算機システムは、DBホスト80、データ位置管理サーバ82、記憶装置10から構成される。DBホスト80、データ位置管理サーバ82、記憶装置10はそれぞれが保有するネットワークインターフェイス78を通してネットワーク79に接続されている。また、DBホスト80、記憶装置10はそれぞれが保有するI/Oバスインターフェイス70からI/Oバス71を介して仮想ボリュームスイッチ72に接続され、これらを通して記憶装置10とDBホスト80間のデータ転送を行う。

10

【0020】

本実施の形態においては、記憶装置10とDBホスト80間のデータ転送を行うI/Oバス71とネットワーク79を異なるものとしているが、例えばiSCSIのような計算機と記憶装置間のデータ転送をネットワーク上で実施する技術も開発されており、本実施の形態においてもこの技術を利用してもよい。このとき、仮想ボリュームスイッチ72において、I/Oバス71とネットワーク79との間でデータ転送が可能であれば記憶装置10もしくはDBホスト80においてI/Oバスインターフェイス70がネットワークインターフェイス78を兼ねても良い。

20

【0021】

記憶装置10は、記憶領域を提供するもので、その記憶領域は記憶領域管理単位であるボリュームを用いて外部に提供し、ボリューム内の部分領域に対するアクセスや管理はブロックを単位として実行する。記憶装置10は、ネットワークインターフェイス78、I/Oバスインターフェイス70、記憶装置制御装置12、ディスクコントローラ16、物理記憶装置18から構成され、ネットワークインターフェイス78、I/Oバスインターフェイス70、記憶装置制御装置12、ディスクコントローラ16はそれぞれ内部バス20により接続され、ディスクコントローラ16と物理記憶装置18は物理記憶装置バス22により接続される。

30

【0022】

記憶装置制御装置12は、CPU24とメモリ26を有する。メモリ26上には、記憶装置におけるキャッシュメモリとして利用するデータキャッシュ28が割り当てられ、記憶装置を制御するためのプログラムである記憶装置制御プログラム40が記憶される。また、メモリ26上には、物理記憶装置18の稼動情報である物理記憶装置稼動情報32、記憶装置10が提供するボリュームを物理的に記憶する物理記憶装置18上の記憶位置の管理情報であるボリューム物理記憶位置管理情報36を保持する。

【0023】

図中の記憶装置10は、複数の物理記憶装置18を有し、1つのボリュームに属するデータを複数の物理記憶装置18に分散配置することが可能である。また、データが記憶される物理記憶装置18上の位置を動的に変更する機能を有する。このデータ移動指示は、ネットワークインターフェイス78を通して外部から行うことが可能である。ただし、それらは必須のものではなく、1つの物理記憶装置の記憶領域をそのままボリュームとして提供する記憶装置10でも本実施の形態にあてはめることができる。

40

【0024】

仮想ボリュームスイッチ72は、記憶装置10から提供される記憶領域管理単位であるボリュームの一部領域または全領域を1つ以上組み合わせた仮想的なボリュームである仮想ボリュームをDBホスト80に提供する機能を有する。仮想ボリュームスイッチ72は、ネットワークインターフェイス78を有し、仮想ボリュームスイッチ72において形成されるボリュームに関する情報である仮想ボリューム情報73を保持する。

【0025】

50

仮想記憶スイッチ 72 はホストから行われたアクセス要求を仮想ボリューム情報 73 を参照して適切な記憶装置 10 の適切な記憶領域へのアクセスへと変換してアクセス処理を実現する。また、ある仮想ボリュームのデータが記憶されるボリュームやそのボリューム内の記憶位置を動的に変更する機能を有する。このデータ移動指示は、ネットワークインターフェイス 78 を通して外部から行うことが可能である。

【 0026 】

DB ホスト 80、データ位置管理サーバ 82 においては、それぞれ CPU 84、ネットワークインターフェイス 78、メモリ 88 を有し、メモリ 88 上にオペレーティングシステム (OS) 100 が記憶・実行されている。

【 0027 】

DB ホスト 80 は I/O パスインターフェイス 70 を有し、仮想ボリュームスイッチ 72 により提供されるボリュームに対してアクセスを実行する。OS 100 内にファイルシステム 104 と 1 つ以上のボリュームからホストが利用する論理的なボリュームである論理ボリュームを作成するボリュームマネージャ 102 と、ファイルシステム 104 やボリュームマネージャ 102 により、OS 100 によりアプリケーションに対して提供されるファイルや論理ローボリュームに記憶されたデータの記録位置等を管理するマッピング情報 106 を有する。

【 0028 】

OS 100 が認識する仮想ボリュームやボリュームマネージャ 102 により提供される論理ボリュームに対して、アプリケーションがそれらのボリュームをファイルと等価なインターフェイスでアクセスするための機構であるロードバース機構を OS 100 が有しているても良い。図中の構成ではボリュームマネージャ 102 が存在しているが、本実施の形態においてはボリュームマネージャ 102 における論理ボリュームの構成を変更することはないので、ボリュームマネージャ 102 が存在せずにファイルシステムが直接仮想ボリュームスイッチ 72 により提供される仮想ボリュームを利用する構成に対しても本実施の形態を当てはめることができる。

【 0029 】

DB ホスト 80 のそれぞれのメモリ 88 上では DBMS 110、データ位置管理副プログラム 120 が記憶・実行され、実行履歴情報 122 が記憶されている。DBMS 110 は内部にスキーマ情報 114 を有している。図中では、DBMS 110 が 1 台のホストに 1 つのみ動作しているが、後述するように、DBMS 110 毎の識別子を用いて管理を行うため、1 台のホストに複数動作していても本実施の形態に当てはめることができる。

【 0030 】

データ位置管理サーバ 82 のメモリ 88 上ではデータ位置管理主プログラム 130 が記憶・実行され、記憶装置 10 内の物理記憶装置 18 の稼動情報である記憶装置稼動情報 132、各記憶装置 10 における物理構成やボリュームの物理記憶位置に関する情報である記憶装置構成情報 134、DB ホスト 80 上の DBMS 110 におけるスキーマに関する情報のうちデータ位置管理に必要なデータである DBMS スキーマ情報 136、DBMS 110 における DBMS 処理の実行履歴に関する情報である DBMS 実行履歴情報 138、DB ホスト 80 における OS 100 内のマッピング情報 106 と仮想ボリュームスイッチ 72 内の仮想ボリューム情報 73 に関する情報を含むデータ記憶位置管理情報 140 が記憶される。

【 0031 】

図中では、データ位置管理サーバ 82 は DB ホスト 80 と異なる計算機となっているが、任意の DB ホスト 80 がデータ位置管理サーバ 82 の役割を兼ねても本実施の形態に当てはめることができる。更に、仮想ボリュームスイッチ 72 上や任意の記憶装置 10 上にデータ位置管理サーバ 82 が提供する機能を実装しても本実施の形態に当てはめることができる。

【 0032 】

図 2 は記憶装置 10 内に保持されている物理記憶装置稼動情報 32 を示す。物理記憶装置

10

20

30

40

50

稼働情報 3 2 中には、記憶装置 1 0 が提供するボリュームの識別子であるボリューム名 5 0 1 とそのボリューム名 5 0 1 を持つボリュームのデータを保持する物理記憶装置 1 8 の識別子である物理記憶装置名 5 0 2、そしてボリューム名 5 0 1 を持つボリュームが物理記憶装置名 5 0 2 を持つ物理記憶装置 1 8 に記憶しているデータをアクセスするための稼働時間のある時刻からの累積値である累積稼働時間 5 0 3 の組を保持する。

【 0 0 3 3 】

記憶装置 1 0 内の記憶装置制御部 1 2 はディスクコントローラ 1 6 を利用して物理記憶装置 1 8 へのデータアクセスする際の開始時刻と終了時刻を取得し、そのアクセスデータがどのボリュームに対するものかを判断して開始時刻と終了時刻の差分を稼働時間として対応するボリューム名 5 0 1 と物理記憶装置名 5 0 2 を持つデータの組の累積稼働時間 5 0 3 に加算する。この情報は、必ずしも全ての記憶装置 1 0 で取得する必要はない。

10

【 0 0 3 4 】

図 3 は記憶装置 1 0 内に保持されているボリューム物理記憶位置管理情報 3 6 を示す。ボリューム物理記憶位置管理情報 3 6 中には、データの論理アドレス-物理記憶装置 1 8 における記憶位置のマッピングを管理するボリューム物理記憶位置メイン情報 5 1 0 と記憶装置 1 0 内でのボリュームに属するデータの物理記憶位置の変更処理の管理情報であるボリュームデータ移動管理情報 5 1 1 が含まれる。

【 0 0 3 5 】

ボリューム物理記憶位置メイン情報 5 1 0 中には、ボリューム名 5 0 1 とそのボリューム上のデータ記憶位置であるボリューム論理ブロック番号 5 1 2 とその論理ブロックが記憶されている物理記憶装置 1 8 の識別子である物理記憶装置名 5 0 2 と物理記憶装置 1 8 上の記憶位置である物理ブロック番号 5 1 4 の組のデータが含まれる。ここで、ボリューム名 5 0 1 が “ E m p t y ” であるエン트리 5 1 5 は特殊なエン트리であり、このエント리는記憶装置 1 0 内の物理記憶装置 1 8 の領域のうち、ボリュームに割り当てられていない領域を示し、この領域に対してデータをコピーすることによりデータの物理記憶位置の動的変更機能を実現する。

20

【 0 0 3 6 】

ボリュームデータ移動管理情報 5 1 1 はボリューム名 5 0 1 と、そのボリューム内の記憶位置を変更するデータ範囲を示す移動論理ブロック番号 7 8 2 と、そのデータが新規に記憶される物理記憶装置 1 8 の識別子とその記憶領域を示す移動先物理記憶装置名 7 8 3 と移動先物理ブロック番号 7 8 4、現在のデータコピー元を示すコピーポインタ 7 8 6 とデータの再コピーの必要性を管理する差分管理情報 7 8 5 の組が含まれる。

30

【 0 0 3 7 】

差分管理情報 7 8 5 とコピーポインタ 7 8 6 を用いたデータの記憶位置変更処理の概略を以下に示す。差分管理情報 7 8 5 はある一定量の領域毎にデータコピーが必要である「 1 」または不必要「 0 」を示すデータを保持する。データの記憶位置変更処理開始時に全ての差分管理情報 7 8 5 のエント리를 1 にセットし、コピーポインタ 7 8 6 を移動元の先頭にセットする。

【 0 0 3 8 】

コピーポインタ 7 8 6 にしたがって差分管理情報 7 8 5 に 1 がセットされている領域を順次移動先にデータをコピーし、コピーポインタ 7 8 6 を更新していく。差分管理情報 7 8 5 で管理される領域をコピーする直前に、その対応するエント리를 0 にセットする。データコピー中に移動領域内のデータに対する更新が行われた場合、それに対応する差分管理情報 7 8 5 のエント리를 1 にセットする。

40

【 0 0 3 9 】

一度全領域のコピーが完了した段階で差分管理情報 7 8 5 内のエント리가全て 0 になったかを確認し、全て 0 であればボリューム物理記憶位置メイン情報 5 1 0 を更新してデータの記憶位置変更処理は完了する。1 のエント리가残っている場合には、再度それに対応する領域をコピーする処理を前記手順で繰り返す。なお、データ記憶位置の動的変更機能の実現方法は他の方式を用いても良い。この場合には、ボリューム物理記憶位置管理情報 3

50

6 中にはボリュームデータ移動管理情報 5 1 1 ではなく他のデータ記憶位置の動的変更機能のための管理情報が含まれることになる。

【 0 0 4 0 】

図 4 は DB ホスト 8 0 の OS 1 0 0 内に記憶されているマッピング情報 1 0 6 を示す。マッピング情報 1 0 6 中には、ボリュームローデバイス情報 5 2 0、ファイル記憶位置情報 5 3 0 と論理ボリューム構成情報 5 4 0 が含まれる。ボリュームローデバイス情報 5 2 0 中には OS 1 0 0 においてローデバイスを指定するための識別子であるローデバイスパス名 5 2 1 とそのローデバイスによりアクセスされる仮想ボリュームあるいは論理ボリュームの識別子であるローデバイスボリューム名 5 2 2 の組が含まれる。

【 0 0 4 1 】

ファイル記憶位置情報 5 3 0 中には、OS 1 0 0 においてファイルを指定するための識別子であるファイルパス名 5 3 1 とそのファイル中のデータ位置を指定するブロック番号であるファイルブロック番号 5 3 2 とそれに対応するデータが記憶されている仮想ボリュームもしくは論理ボリュームの識別子であるファイル配置ボリューム名 5 3 3 とそのボリューム上のデータ記憶位置であるファイル配置ボリュームブロック番号 5 3 4 の組が含まれる。

【 0 0 4 2 】

論理ボリューム構成情報 5 4 0 中にはボリュームマネージャ 1 0 2 により提供される論理ボリュームの識別子である論理ボリューム名 5 4 1 とその論理ボリューム上のデータ的位置を示す論理ボリューム論理ブロック番号 5 4 2 とその論理ブロックが記憶されている仮想ボリュームの識別子である仮想ボリューム名 5 4 3 と仮想ボリューム上の記憶位置である仮想ボリュームブロック番号 5 4 4 の組が含まれる。

【 0 0 4 3 】

図 5 は DBMS 1 1 0 内に記憶されているその内部で定義・管理しているデータその他の管理情報であるスキーマ情報 1 1 4 を示す。スキーマ情報 1 1 4 には、表のデータ構造や制約条件等の定義情報を保持する表定義情報 5 5 1、索引のデータ構造や対象である表等の定義情報を保持する索引定義情報 5 5 2、利用するログに関する情報であるログ情報 5 5 3、利用する一時表領域に関する情報である一時表領域情報 5 5 4、管理しているデータのデータ記憶位置の管理情報であるデータ記憶位置情報 5 5 5、キャッシュの構成に関する情報であるキャッシュ構成情報 5 5 6 とデータをアクセスする際の並列度に関する情報である最大アクセス並列度情報 5 5 7 を含む。

【 0 0 4 4 】

データ記憶位置情報 5 5 5 中には、表、索引、ログ、一時表領域等のデータ構造の識別子であるデータ構造名 5 6 1 とそのデータを記憶するファイルまたはローデバイスの識別子であるデータファイルパス名 5 6 2 とその中の記憶位置であるファイルブロック番号 5 6 3 の組が含まれる。

【 0 0 4 5 】

キャッシュ構成情報 5 5 6 は DBMS 1 1 0 が三種類のキャッシュ管理のグループを定義し、そのグループに対してキャッシュを割り当てている場合を示す。キャッシュ構成情報 5 5 6 中には、グループ名 5 6 5 とグループ中のデータ構造のデータをホスト上にキャッシュする際の最大データサイズであるキャッシュサイズ 5 6 6 とそのグループに所属するデータ構造の識別子の所属データ構造名 5 6 7 の組が含まれる。最大アクセス並列度情報 5 5 7 には、データ構造名 5 6 1 とそのデータ構造にアクセスする際の一般的な場合の最大並列度に関する情報である最大アクセス並列度 5 6 9 の組が含まれる。

【 0 0 4 6 】

図 6 は DB ホスト 8 0 のメモリ 8 8 上に記憶されている実行履歴情報 1 2 2 を示す。実行履歴情報 1 2 2 中には、DBMS 1 1 0 で実行されたクエリ 5 7 0 の履歴が記憶されている。この情報は、DBMS 1 1 0 が作成する。または DBMS のフロントエンドプログラムがこの情報を作成する。この場合には、DBMS フロントエンドプログラムが存在する計算機に実行履歴情報 1 2 2 が記憶されることになる。

10

20

30

40

50

【 0 0 4 7 】

図 7 は仮想ボリュームスイッチ 7 2 が保持する仮想ボリューム情報 7 3 を示す。仮想ボリューム情報 7 3 は仮想ボリューム記憶位置情報 7 9 0 と仮想ボリュームデータ移動管理情報 7 9 1 を含む。仮想ボリューム記憶位置情報 7 9 0 中には、仮想ボリュームスイッチ 7 2 が DB ホスト 8 0 に提供する仮想ボリュームの識別子である仮想ボリューム名 5 4 3 とその仮想ボリューム上のデータの記憶位置を示す仮想ボリュームブロック番号 5 4 4 とそのブロックが記憶されている記憶装置 1 0 の識別子である記憶装置名 5 8 3 とそのボリュームの識別子であるボリューム名 5 1 1 とボリューム上の記憶位置であるボリューム論理ブロック番号 5 1 2 の組が含まれる。

【 0 0 4 8 】

仮想ボリューム名 5 4 3 が “ E m p t y ” であるエントリ 5 8 5 は特殊なエントリであり、このエントリに含まれる記憶装置 1 0 上の領域は DB ホスト 8 0 に対して仮想ボリュームとして提供されていない領域を示す。これらの領域に対して仮想ボリュームスイッチ 7 2 はデータの移動を行うことができる。

【 0 0 4 9 】

仮想ボリュームデータ移動管理情報 7 9 1 は仮想ボリューム名 5 4 3 と、そのボリューム内の記憶位置を変更するデータ範囲を示す移動仮想ボリュームブロックブロック番号 7 9 3 と、そのデータが新規に記憶される記憶装置 1 0 の識別子とその記憶領域を示す移動先記憶装置名 7 9 4 と移動先ボリューム名 7 9 5 と移動先論理ブロック番号 7 9 6、現在のデータコピー元を示すコピーポインタ 7 8 6 とデータの再コピーの必要性を管理する差分管理情報 7 8 5 の組が含まれる。

【 0 0 5 0 】

図 3 中のボリュームデータ移動管理情報 5 1 1 で説明したのと同様の方式によりデータの記憶位置の動的変更機能を実現できる。データ記憶位置の動的変更機能の実現方法は他の方式を用いても良い。この場合には、仮想ボリューム情報 7 3 中には仮想ボリュームデータ移動管理情報 7 9 1 ではなく他のデータ記憶位置の動的変更機能のための管理情報が含まれることになる。

【 0 0 5 1 】

図 8 はデータ位置管理サーバ 8 2 上に記憶される記憶装置稼働情報 1 3 2 を示す。記憶装置稼働情報 1 3 2 中には、記憶装置 1 0 の識別子である記憶装置名 5 8 3、記憶装置 1 0 が提供するボリュームの識別子であるボリューム名 5 0 1、記憶装置 1 0 に存在する物理記憶装置 1 8 の識別子である物理記憶装置名 5 0 2、記憶装置名 5 8 3 とボリューム名 5 0 1、物理記憶装置名 5 0 2 により特定される領域の前回取得時の累積稼働時間 5 0 3 の値である旧累積稼働時間 5 9 3 とある一定時間内の動作時間の割合を示す稼働率 5 9 4 の組と、稼働率計算のために前回累積稼働時間を取得した時刻である前回累積稼働時間取得時刻 5 9 5 を含む。

【 0 0 5 2 】

記憶装置 1 0 は物理記憶装置稼働情報 3 2 を外部に提供する機構を有し、それを利用して記憶位置管理主プログラム 1 3 0 は記憶装置 1 0 で取得・記憶されている物理記憶装置稼働情報 3 2 をネットワーク 7 9 を通して一定間隔で取得し、取得した累積稼働時間 5 0 3 と旧累積稼働時間 5 9 3、前回累積稼働時間取得時刻 5 9 5 と現データ取得時刻を用いて前回累積稼働時間取得時刻 5 9 5 と現データ取得時刻間の稼働率 5 9 4 を計算・記憶する。その後、取得した累積稼働時間 5 0 3 を旧累積稼働時間 5 9 3 に、現データ取得時刻を前回累積稼働時間取得時刻 5 9 5 に記憶する。

【 0 0 5 3 】

なお、全ての記憶装置 1 0 で物理記憶装置稼働情報 3 2 を保持しているとは限らない。その場合には物理記憶装置稼働情報 3 2 を保持している記憶装置 1 0 についてのみ記憶装置稼働情報 1 3 2 のエントリに含める。また、全ての記憶装置 1 0 で物理記憶装置稼働情報 3 2 を保持していない場合には記憶装置稼働情報 1 3 2 を保持しなくてもよい。

【 0 0 5 4 】

10

20

30

40

50

図9はデータ位置管理サーバ82上に記憶される記憶装置構成情報134を示す。記憶装置構成情報134中には、記憶装置10の識別子である記憶装置名583と、記憶装置10がデータ記憶位置の動的変更機能を有しているかいないかの情報である移動機能情報601と記憶装置10が保持しているデータキャッシュの容量であるデータキャッシュ容量602、記憶装置名583を持つ記憶装置10におけるボリューム物理記憶位置メイン情報510を保持する記憶装置ボリューム物理記憶位置管理情報603の組を保持する。

【0055】

記憶装置10はボリューム物理記憶位置メイン情報510とデータキャッシュ28のサイズに関する情報を外部に提供する機能を有し、記憶装置構成情報134を作成するため、データ記憶位置管理主プログラム130は記憶装置10からボリューム物理記憶位置情報36とデータキャッシュ28のサイズに関する情報をネットワーク79を通して取得する。記憶装置10はデータキャッシュ28のサイズは必ずしも外部に提供する機能を有する必要はなく、その場合には、データキャッシュ容量602の対応部分はデータ無効を記憶しておく。

【0056】

一方、ボリューム物理記憶位置メイン情報510に関しては、記憶装置10が外部に提供する機能を有さなくても構わない場合は、物理記憶装置18を1つしか保持せずにそれをそのまま1つのボリュームとして提供する等、記憶装置10が提供するボリュームがどのように物理記憶装置18上に記憶されるかが固定され、かつ、そのマッピングをあらかじめデータ記憶位置管理主プログラム130が理解している場合である。このとき、データ記憶位置管理主プログラム130はこのルールに従って記憶装置ボリューム物理記憶位置管理情報603の内容を設定する。このルールは設定ファイル等を用いて管理者がデータ記憶位置管理主プログラム130に与える。

【0057】

図10はデータ位置管理サーバ82上に記憶されるDBMSスキーマ情報136を示す。DBMSスキーマ情報136は、DBMSデータ構造情報621、DBMSデータ記憶位置情報622、DBMSパーティション化表・索引情報623、DBMS索引定義情報624、DBMSキャッシュ構成情報625、DBMSホスト情報626を含む。

【0058】

DBMSデータ構造情報621はDBMS110で定義されているデータ構造に関する情報で、DBMS110の識別子であるDBMS名631、DBMS110内の表・索引・ログ・一時表領域等のデータ構造の識別子であるデータ構造名561、データ構造の種類を表すデータ構造種別640、そのデータ構造をアクセスする際の最大並列度に関する情報である最大アクセス並列度569の組を保持する。このとき、データ構造によっては最大アクセス並列度569の値を持たない。

【0059】

DBMSデータ記憶位置情報622はDBMS名631とそのDBMSにおけるデータ記憶位置管理情報555であるデータ記憶位置管理情報638の組を保持する。DBMSパーティション化表・索引情報623は、1つの表や索引をある属性値により幾つかのグループに分割したデータ構造を管理する情報で、パーティション化されたデータ構造が所属するDBMS110の識別子であるDBMS名631と分割化される前のデータ構造の識別子であるパーティション元データ構造名643と分割後のデータ構造の識別子であるデータ構造名561とその分割条件を保持するパーティション化方法644の組を保持する。今後、パーティション化されたデータ構造に関しては、特に断らない限り単純にデータ構造と呼ぶ場合にはパーティション化後のものを指すものとする。

【0060】

DBMS索引定義情報624には、DBMS名631、索引の識別子である索引名635、その索引のデータ形式を示す索引タイプ636、その索引がどの表のどの属性に対するものかを示す対応表情報637の組を保持する。DBMSキャッシュ構成情報625は、DBMS110のキャッシュに関する情報であり、DBMS名631とDBMS110に

10

20

30

40

50

おけるキャッシュ構成情報 5 5 6 の組を保持する。DBMS ホスト情報 6 2 6 は、DBMS 名 6 3 1 を持つ DBMS 1 1 0 がどのホスト上で実行されているかを管理するもので、DBMS 名 6 3 1 と DBMS 実行ホストの識別子であるホスト名 6 5 1 の組を保持する。

【 0 0 6 1 】

DBMS スキーマ情報 1 3 6 中の DBMS ホスト情報 6 2 6 以外は、データ記憶位置管理主プログラム 1 3 0 が DBMS 1 1 0 が管理しているスキーマ情報 1 1 4 の中から必要な情報を取得して作成するものである。DBMS 1 0 0 のスキーマ情報 1 1 4 は、ネットワーク 7 9 を通してデータ位置管理主プログラム 1 3 0 が直接、あるいは、データ位置管理副プログラム 1 2 0 を介して、SQL 等のデータ検索言語を用いてビューとして公開されている情報を取得するか、または、専用の機構を用いて取得する。DBMS ホスト情報 6 2 6 はシステム構成情報であり、管理者が設定する。

10

【 0 0 6 2 】

図 1 1 はデータ位置管理サーバ 8 2 上に記憶される DBMS 実行履歴情報 1 3 8 を示す。DBMS 実行履歴情報 1 3 8 には、DBMS 1 1 0 の識別子の DBMS 名 6 3 1 と各々の DBMS 1 1 0 で実行されたクエリ 5 7 0 の履歴が保持される。これは、データ位置管理主プログラム 1 3 0 がネットワーク 7 9 を通して DB ホスト 8 0 内に記憶されている実行履歴情報 1 2 2 をデータ位置管理副プログラム 1 2 0 を利用して収集し、それを記憶したものである。

【 0 0 6 3 】

また、前述のように実行履歴情報 1 2 2 が DBMS フロントエンドプログラムが実行される計算機上に記憶される可能性も存在する。この場合には、DBMS フロントエンドプログラムが実行される計算機から実行履歴情報 1 2 2 をデータ位置管理サーバ 8 2 へ転送する手段を設け、転送された実行履歴情報 1 2 2 をデータ位置管理主プログラム 1 3 0 が DBMS 実行履歴情報 1 3 8 として記憶する。なお、本実施の形態においては、全ての DBMS 1 1 0 から実行履歴情報 1 2 2 を収集する必要はなく、DBMS 実行履歴情報 1 3 8 は存在しなくてもよい。

20

【 0 0 6 4 】

図 1 2 はデータ位置管理サーバ 8 2 上に記憶されるデータ記憶位置管理情報 1 4 0 を示す。データ記憶位置管理情報 1 4 0 には、ホストマッピング情報 6 5 0 と仮想ボリューム記憶位置管理情報 7 9 0 が含まれる。ホストマッピング情報 6 5 0 には、ホストの識別子であるホスト名 6 5 1 とそのホストにおけるマッピング情報 1 0 6 の組が保持される。これは、データ位置管理主プログラム 1 3 0 がネットワーク 7 9 を通して、DB ホスト 8 0 の OS 1 0 0 が保持しているマッピング情報 1 0 6 を、データ位置管理副プログラム 1 2 0 を利用して収集し、それを記憶したものである。

30

【 0 0 6 5 】

データ位置管理副プログラム 1 2 0 は、OS 1 0 0 が提供している管理コマンドや情報提供機構、参照可能な管理データの直接解析等によりマッピング情報 1 0 6 を取得する。仮想ボリュームスイッチ 7 2 は外部に仮想ボリューム記憶位置管理情報 7 9 0 を提供する機構を有し、データ位置管理主プログラム 1 3 0 はネットワーク 7 9 を介して仮想ボリュームスイッチ 7 2 から仮想ボリューム記憶位置管理情報 7 9 0 を取得する。

40

【 0 0 6 6 】

図 1 3 はデータ位置管理主プログラム 1 3 0 によるデータ再配置処理の処理フローを示す。ここで、処理開始は管理者の指示によることとする。後述するように、複数の異なった種類のデータ配置解析・データ再配置案作成処理を実行可能であり、処理すべき種類の指定をして処理を開始する。また、処理にパラメータが必要な場合は併せてそれを管理者が指示をするものとする。本実施の形態においては、データの記憶位置の動的変更機能は仮想ボリュームスイッチ 7 2 と記憶装置 1 0 が保持する。ただし、記憶装置 1 0 はデータの記憶位置の動的変更機能を有さなくてもよい。

【 0 0 6 7 】

ステップ 2 0 0 1 でデータ再配置処理を開始する。このとき、データ配置解析・データ再


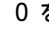
50

配置案作成処理として何を実行するか指定する。また、必要であればパラメータを指定する。

【0068】

ステップ2002でデータ再配置処理に必要な情報を収集し、記憶装置稼働情報132、記憶装置構成情報134、DBMSスキーマ情報136、DBMS実行履歴情報138、データ記憶位置管理情報140として記憶する。なお、このデータ収集は、ステップ2001の処理開始とは無関係にあらかじめ実行しておくこともできる。この場合には、情報を取得した時点から現在まで情報に変更がないかどうかをこのステップで確認する。

【0069】

ステップ2003では、ワーク領域を確保し、その初期化を行う。ワーク領域としては、14に示すデータ再配置ワーク情報670と作成した15に示す移動プラン情報750を利用する。データ再配置ワーク情報670と移動プラン情報750の詳細とその初期データ作成方法は後述する。

10

【0070】

ステップ2004でデータ配置の解析・再配置案の作成処理を実行する。後述するように、データ配置の解析・再配置案作成処理は複数の観点による異なったものが存在し、このステップではステップ2001で指定された処理を実行する。またステップ2001でパラメータを受け取った場合には、それを実行する処理に与える。

【0071】

ステップ2005ではステップ2004のデータ再配置案作成処理が成功したかどうか確認する。成功した場合にはステップ2007に進む。失敗した場合にはステップ2006に進み、管理者にデータ再配置案作成が失敗したことを通知し、ステップ2010に進み処理を完了する。

20

【0072】

ステップ2007では、ステップ2004で作成されたデータ再配置案を管理者に提示する。この提示を受けて管理者はデータ再配置案に問題がないか判断する。

ステップ2008では、データの再配置を続行するか否かを管理者から指示を受ける。続行する場合にはステップ2009に進み、そうでない場合にはステップ2010に進み処理を完了する。

【0073】

ステップ2009では、ステップ2004で作成されたデータの再配置案を基にデータの再配置指示を仮想ボリュームスイッチ72、あるいは、記憶装置10に対して発行する。仮想ボリュームスイッチ72と記憶装置10はネットワーク79を通してデータのデータ再配置指示を受ける機能を有し、これを利用して指示が出される。

30

【0074】

この指示の形式としては、仮想ボリュームスイッチ72に対しては、それが提供する仮想ボリュームのあるデータ領域を、記憶装置10のボリューム内のデータ領域を指定してそこへの移動を指示するものとなり、記憶装置10に対しては、それが提供するボリュームのあるデータ領域を、その記憶装置10内の物理記憶装置18のデータ領域を指定してそこへの移動を指示するものとなる。この指示に従って、仮想ボリュームスイッチ72、あるいは、記憶装置10はデータの再配置処理を実行する。

40


【0075】

ステップ2010でデータ再配置処理は完了である。

【0076】

この処理フローは、管理者の指示で処理を開始し、作成されたデータ再配置案に問題点がないかステップ2007と2008で管理者が判断している。この管理者による確認を省略し、タイマを用いて設定した処理開始時刻にデータ配置解析・再配置案作成処理を開始することによりデータ再配置処理の自動化も可能である。

【0077】

14は13のステップ2003において作成する情報であるデータ再配置ワーク情報

50

670を示す。データ再配置ワーク情報670中には、仮想ボリューム物理記憶位置情報680とデータ構造仮想ボリューム内位置情報690が含まれる。

【0078】

仮想ボリューム物理記憶位置情報680は仮想ボリュームスイッチ72が提供する仮想ボリュームのデータがどの記憶装置10内のどの物理記憶装置18のどの位置に記憶されているかに関する情報であり、仮想ボリューム名543、仮想ボリュームブロック番号544、記憶装置名583、ボリューム名501、物理記憶装置名502、物理ブロック番号514の組で表される。仮想ボリューム物理記憶位置情報680では、データ記憶位置管理情報140中の仮想ボリューム記憶位置管理情報と、記憶装置構成情報134中の記憶装置名583と記憶装置ボリューム物理記憶位置情報603を参照し、記憶装置10が提供

10

【0079】

仮想ボリューム名543が“Empty”であるエントリ群681は記憶装置10内の物理記憶装置18の記憶領域のうち、データ再配置のためにデータを移動することが可能な領域の集合を表し、データ再配置案作成時に、この領域の中から適切なデータの移動先を発見する。このうち、ボリューム名501が有効値のエントリは、仮想ボリュームスイッチ72におけるデータ記憶位置の動的変更に利用可能な領域であり、その利用には特に制約はない。一方のボリューム名501が無効値のエントリは、記憶装置10におけるデータ記憶位置の動的変更に利用可能な領域であり、その領域が属する記憶装置10内に記憶

20

【0080】

データ構造仮想ボリューム内位置情報690は、DBMS110が保持しているデータ構造が仮想ボリュームスイッチ72が提供する仮想ボリュームのどこに記憶されているかを示した情報であり、DBMS名631、データ構造名561、データファイルパス名562、ファイルブロック番号563、仮想ボリューム名543、仮想ボリュームブロック番号544の組を保持する。この情報では、DBMSスキーマ情報136中のDBMSデータ記憶位置情報622とDBMSホスト情報626とデータ記憶位置管理情報140中のホストマッピング情報650を参照し、ファイル(ローデバイス)パス、論理(仮想)ボリュームに関して対応する部分をまとめることにより初期データが作成される。

【0081】

図15は図13のステップ2004で実行されるデータ配置解析・データ再配置案作成処理により作成されるデータ移動案を格納する移動プラン情報750を示す。移動プラン情報750は、仮想ボリュームスイッチ72に対するデータ移動指示を記憶する仮想ボリューム移動プラン情報751と記憶装置10に対するデータ移動指示を記憶する物理記憶位置移動プラン情報752を含む。これら情報に関しては、何もデータを持たないように初期化する。

30

【0082】

仮想ボリューム移動プラン情報751には、移動指示の実行順序を示す移動順序761、移動するデータを持つ仮想ボリュームとそのデータ領域を示す移動仮想ボリューム名762と移動仮想ボリュームブロック番号763、そのデータの移動先の記憶装置、ボリューム、ボリューム内の記憶領域を示す移動先記憶装置名764と移動先ボリューム名765と移動先ボリューム論理ブロック番号766の組を保持する。

40

【0083】

物理記憶位置移動プラン情報752には、移動指示の実行順序を示す移動順序761、移動するデータを持つ記憶装置10とボリュームとそのデータ領域を示す移動記憶装置名767と移動ボリューム名768と移動ボリューム論理ブロック番号769、そのデータの移動先の物理記憶装置とその記憶領域を示す移動先物理記憶装置名771と移動先物理ブロック番号772の組を保持する。なお、物理記憶位置移動プラン情報752に関しては、全ての記憶装置10が記憶装置内の物理記憶位置の動的変更機能を有さない場合には、この情報を保持する必要はない。

50

【 0 0 8 4 】

図 1 3 のステップ 2 0 0 4 で実行されるデータ配置解析・データ再配置案作成処理について説明する。前述のように、この処理には幾つかの種類が存在する。全ての処理に共通するのは逐次的にデータ再配置するためのデータ移動案を作成することである。そのため、データ移動の順番には意味があり、移動プラン情報 7 5 0 中の移動順序 7 6 1 にその順番を保持し、その順序どおりにデータ移動を行うことによりデータの再配置を実施する。また、逐次処理のため、移動後のデータ配置をもとに次のデータの移動方法を決定する必要がある。そこで、データ移動案を作成するたびにデータ再配置ワーク情報 6 7 0 をデータ移動後の配置に更新する。

【 0 0 8 5 】

データ再配置案作成時のデータ移動案の作成は以下のように行う。移動したいデータ量以上の連続した移動可能領域をデータ再配置ワーク情報 6 7 0 中の情報から抽出し、その中の領域を適当に選択し、設定条件や後述する制約を満たすかどうか確認をする。もし、それらを満たす場合にはそこを移動先として設定する。それらを満たさない場合には他の領域を選択し、再度それらを満たすかどうか確認をする。以下、設定条件と制約を満たす領域を検出するか、全ての移動したいデータ量以上の連続した移動可能領域が設定条件や制約を満たさないことを確認するまで処理を繰り返す。もし、全ての領域で設定条件や制約を満たさない場合にはデータ移動案の作成に失敗として終了する。

【 0 0 8 6 】

このときに重要なのは移動後のデータ配置において、問題となる配置を行わないことである。特に R D B M S においては、特定のデータに関してはアクセスが同時に行われる可能性が高く、それらを異なる物理記憶装置 1 8 上に配置する必要がある。

【 0 0 8 7 】

そこで、以下で説明する全てのデータの移動案を作成する場合には、移動するデータに含まれるデータ構造と、移動先に含まれるデータ構造を調べ、ログとその他のデータ、一時表領域とその他のデータ、表データとそれに対して作成された木構造の索引データがデータの移動後に同じ物理記憶装置 1 8 に配置されるかどうかを確認し、配置される場合には、その配置案は利用不可能と判断する。

【 0 0 8 8 】

なお、あるデータ構造がどの記憶装置 1 0 のどのボリューム上の領域に、あるいは、どの物理記憶装置 1 8 の領域に記憶されているか、また逆に、ある物理記憶装置 1 8 上や記憶装置 1 0 のボリュームの領域に記憶されるデータがどのデータ構造に対応するかは、データ再配置ワーク情報 6 7 0 中の仮想ボリューム物理記憶位置情報 6 8 0 とデータ構造仮想ボリューム内位置情報 6 9 0 を仮想ボリュームに関して対応する部分を組み合わせることにより把握可能である。

【 0 0 8 9 】

図 1 6 に第 1 のデータ配置解析・データ再配置案作成処理である、記憶装置稼働情報 1 3 2 を基にした同時アクセス実行データ構造を分離するためのデータ再配置案作成処理の処理フローを示す。本処理においては、物理記憶装置 1 8 の稼働率が閾値を超えたものはディスクネック状態にあると判断してそれを解消するデータの移動案を作成する。

【 0 0 9 0 】

前述のように、記憶装置稼働情報 1 3 2 は全ての記憶装置 1 0 中の物理記憶装置 1 8 に関する稼働情報を必ずしも含む訳ではない。稼働情報が存在しない物理記憶装置 1 8 に関しては、本処理におけるデータ再配置案作成の対象外とし、それらが存在しないものとして扱う。本処理は、実測値に基づいて問題点を把握し、それを解決する方法を見つけるため、より精度の高いデータ再配置案を作成するものであり、ボトルネックの自動解決方式としてデータ移動の自動化機能に組み入れても効果的に働く。

【 0 0 9 1 】

ステップ 2 1 0 1 で処理を開始する。本処理を開始するにあたっては、どの期間の稼働率を参照するかを指定する。

10

20

30

40

50

【 0 0 9 2 】

ステップ 2 1 0 2 では、物理記憶装置 1 8 の識別子と指定期間における物理記憶装置 1 8 の稼働率の組を記憶するワーク領域を取得し、記憶装置稼働情報 1 3 2 を参照してその情報を設定し、物理記憶装置 1 8 の稼働率で降順にソートする。記憶装置稼働情報 1 3 2 中では、同じ物理記憶装置 1 8 中に記憶されているデータであっても異なるボリュームのものは分離して稼働率を取得しているため、それらの総和として物理記憶装置 1 8 の稼働率を求める必要がある。ステップ 2 1 0 3 では、ステップ 2 1 0 2 のソート結果をもとに物理記憶装置 1 8 の稼働率が閾値を超えているもののリストである過負荷確認リストを作成する。このリスト中のエントリに関して稼働率が降順になるような順序を保つようにする。

10

【 0 0 9 3 】

ステップ 2 1 0 4 では、過負荷確認リスト中にエントリが存在するか確認する。エントリが存在しない場合には、もう過負荷状態の物理記憶装置 1 8 が存在しないものとしてステップ 2 1 0 5 に進みデータ再配置案作成処理成功として処理を終了する。エントリが存在する場合には、ステップ 2 1 0 6 に進む。

【 0 0 9 4 】

ステップ 2 1 0 6 では、過負荷確認リスト中の最も物理記憶装置 1 8 の稼働率が高いものを再配置対象の物理記憶装置 1 8 として選択する。ステップ 2 1 0 7 では、再配置対象となった物理記憶装置 1 8 内部のボリュームとその稼働率のリストを記憶装置稼働情報 1 3 2 を参照して作成し、稼働率で降順にソートする。

20

【 0 0 9 5 】

ステップ 2 1 0 8 では、リスト中のあるボリュームの稼働率があらかじめ定められた閾値を超過しているかどうか確認する。全てのボリュームの稼働率が閾値を超えていない場合には、ステップ 2 1 1 3 に進み、あるボリュームの稼働率がその閾値を超えている場合には、ステップ 2 1 0 9 に進む。

【 0 0 9 6 】

ステップ 2 1 0 9 においては、稼働率が閾値を超えているボリュームに関して、確認対象の物理記憶装置 1 8 中に同時にアクセスされる可能性があるデータの組、すなわち、ログとその他のデータ、一時表領域とその他のデータ、表データとそれに対して作成された木構造の索引データがあるそのボリューム内部に記憶されているかどうかを検出する処理を行う。ステップ 2 1 1 0 では、ステップ 2 1 0 9 における結果を確認し、同時アクセスデータ構造の組が存在する場合にはステップ 2 1 1 1 に進む。同時アクセスデータ構造の組が存在しない場合には、ステップ 2 1 1 2 に進む。

30

【 0 0 9 7 】

ステップ 2 1 1 1 においては、同時アクセスデータ構造の組に属するデータを異なる物理記憶装置 1 8 に記憶するためのデータ移動案を作成し、ステップ 2 1 1 4 に進む。

【 0 0 9 8 】

ステップ 2 1 1 2 においては、現在確認対象となっているのボリューム内のデータを論理ブロック番号に従って 2 分割し、その片方を他の物理記憶装置 1 8 へ移動するデータ移動案を作成し、ステップ 2 1 1 4 に進む。

40

【 0 0 9 9 】

ステップ 2 1 1 3 においては、現在確認対象になっている物理記憶装置 1 8 の稼働率が閾値を下回るまで、稼働率が高いボリュームから順に、その物理記憶装置 1 8 に記憶されているボリュームを構成するデータ全体を他の物理記憶装置 1 8 に移動するデータ移動案を作成し、ステップ 2 1 1 4 に進む。

【 0 1 0 0 】

ステップ 2 1 1 1 , 2 1 1 2 , 2 1 1 3 のデータ移動先を検出する際に、移動後の移動先の記憶装置の稼働率を予測する。物理記憶装置 1 8 毎の性能差が既知の場合にはその補正を行った移動データを含む記憶装置 1 8 上のボリュームの稼働率分、未知の場合には補正を行わない移動データを含む記憶装置 1 8 上のボリュームの稼働率分、データ移動により

50

移動先の物理記憶装置 18 の稼働率が上昇すると考え、加算後の値が閾値を越えないような場所へのデータの移動案を作成する。稼働率の加算分に関して、移動データ量の比率を考慮しても良いが、ここではデータ中のアクセスの偏りを考慮して移動データに全てのアクセスが集中したと考えた判断を行う。

【0101】

ステップ 2114 では、データの移動案の作成に成功したかどうかを確認し、失敗した場合にはステップ 2117 に進みデータの再配置案作成処理失敗として処理を終了する。成功した場合にはステップ 2115 に進む。

【0102】

ステップ 2115 では作成したデータ移動案を移動プラン情報 750 に追加、ステップ 2116 に進む。ステップ 2116 ではデータ再配置ワーク情報 670 を作成したデータ移動案に従って修正し、移動先記憶装置 18 のステップ 2102 で作成した物理記憶装置 18 毎の稼働情報情報の値を前述の移動後の稼働率判断値に修正する。その後、現在の確認対象の物理記憶装置 18 を過負荷確認リストから削除し、ステップ 2104 に戻り次の確認を行う。

【0103】

次に第 2 のデータ配置解析・データ再配置案作成処理である、DBMS 実行履歴情報 138 を基にした同時アクセス実行データ構造を分離するためのデータ再配置案作成処理を示す。本処理においては、クエリの実行履歴から同時にアクセスされるデータの組を取得し、それらを異なる物理記憶装置 18 に配置するデータ再配置案を作成する。前述のように、全ての DBMS 110 について実行履歴を取得できるわけではない。本処理実行時に実行履歴が存在しない DBMS 110 が利用するデータに関してはデータ再配置の対象外とする。

【0104】

本処理においては、図 17 に示すクエリ実行時同時アクセスデータカウント情報 700 を利用する。クエリ実行時同時アクセスデータカウント情報 700 は、DBMS 名 631、同時にアクセスされる可能性のあるデータ構造のデータ構造名 561 の組を示すデータ構造名 A 701 とデータ構造名 B 702、そして、DBMS 実行履歴情報 138 の解析によりそのデータ構造の組がアクセスされたと判断された回数であるカウント値 703 の組で表される。この組はカウント値 703 の値でソートしておく。

【0105】

クエリ実行時同時アクセスデータカウント情報 700 は DBMS 実行履歴情報 138 から作成する。最初にクエリ実行時同時アクセスデータカウント情報 700 のエントリを全消去する。DBMS 100 において定型処理が行われる場合には、まず、その型により分類し、その型の処理が何回実行されたかを確認する。

【0106】

続いて DBMS 100 から型毎のクエリ実行プランを取得する。そのクエリ実行プランにより示される処理手順から同時にアクセスされるデータ構造の組を判別する。そして、クエリ実行時同時アクセスデータカウント情報 700 中の DBMS 名 631・データ構造名 A 701・データ構造名 B 702 を参照し、既に対応するデータ構造の組が存在している場合には先に求めたその型の処理回数をカウント値 703 に加算する。既に対応するデータ構造の組が存在していない場合には、新たにエントリを追加してカウント値 703 を先に求めたその型の処理回数にセットする。

【0107】

DBMS 100 において非定型処理が行われる場合には、1つ1つの実行されたクエリに関してクエリ実行プランを取得し、そのクエリ実行プランにより示される処理手順から同時にアクセスされるデータ構造の組を判別する。そして、クエリ実行時同時アクセスデータカウント情報 700 中の DBMS 名 631・データ構造名 A 701・データ構造名 B 702 を参照し、既に対応するデータ構造の組が存在している場合にはカウント値 703 に 1 を加算する。既に対応するデータ構造の組が存在していない場合には、新たにエントリ

10

20

30

40

50

を追加してカウント値 703 に 1 をセットする。

【0108】

クエリ実行プランから同時にアクセスされる可能性があるデータ構造の判別は以下のように行う。まず、木構造の索引に対するアクセスが実施される場合には、その木構造索引データと、その索引が対象とする表データが同時にアクセスされると判断する。また、データの更新処理や挿入処理が行われる場合には、ログとその他のデータが同時にアクセスされると判断する。以下は DBMS 110 の特性に依存するが、例えば、クエリ実行プラン作成時にネストループジョイン処理を多段に渡り実行する計画を作成し、それらの多段に渡る処理を同時に実行する RDBMS が存在する。この RDBMS を利用する場合にはその多段に渡るネストループジョイン処理で利用する表データとその表に対する木構造の索引データは同時にアクセスされると判断できる。

10

【0109】

このように、クエリ実行計画による同時アクセスデータの判断に関しては、DBMS 110 の処理特性を把握して判断する必要があるが、ここでは、対象とする DBMS 110 の種類を絞りデータ位置管理主プログラム 130 が DBMS 110 特有の同時アクセスデータ構造の組を把握できる機能を有することを仮定する。

【0110】

図 18 に DBMS 実行履歴情報 138 を基にした同時アクセス実行データ構造を分離するためのデータ再配置案作成処理処理フローを示す。ステップ 2201 で処理を開始する。ステップ 2202 で実行履歴から同時にアクセスされるデータ構造の組とその実行頻度に関する情報である、前述のクエリ実行同時アクセスデータ構造カウント情報 700 を作成する。ステップ 2203 において、カウント値 703 の全エントリの総和に対してカウント値 703 の値が一定割合以上のデータ構造とその所属する DBMS 110 の組を求め、それらを確認リストとして記憶する。

20

【0111】

ステップ 2204 でステップ 2203 で求めた確認リスト中に含まれるデータ構造の組に関して、それらを異なる物理記憶装置 18 に記憶するデータ再配置案を作成し、ステップ 2205 に進む。なお、ステップ 2204 の処理に関しては、図 19 を用いて後で説明する。ステップ 2205 では、ステップ 2204 においてデータ再配置案の作成に成功したかどうかを確認し、成功した場合にはステップ 2206 に進みデータ再配置案作成処理成功として処理を終了し、失敗した場合にはステップ 2207 に進みデータ再配置案作成処理失敗として処理を終了する。

30

【0112】

図 19 に、指定されたデータ構造とそのデータ構造と同時にアクセスされる可能性が高いデータ構造の組を分離するデータ再配置案を作成する処理のフローを示す。本処理を開始するときには、データ構造名と物理記憶装置 18 から分離するデータ構造名の組のリストである確認リストを与える。

【0113】

ステップ 2301 で処理を開始する。ステップ 2303 で確認リスト中にエントリが存在するか確認し、存在しない場合にはステップ 2304 に進みデータ再配置案作成処理成功として処理を終了する。存在する場合にはステップ 2305 に進む。

40

【0114】

ステップ 2305 においては、確認リストから 1 つ確認対象データ構造名とその所属 DBMS 名の組とその分離データ構造名とその所属 DBMS 名の組の組を取得し、ステップ 2306 に進む。

【0115】

ステップ 2306 においては、確認対象データ構造とその分離するデータ構造が同一の物理記憶装置上に記憶されているかどうかの確認を行う。前述のように、この確認はデータ再配置ワーク情報 670 を参照することにより可能である。両データ構造が全て異なる物理記憶装置上に存在する場合にはステップ 2312 に進み、ある物理記憶装置上に両デー

50

タ構造が存在する場合にはステップ2307に進む。

【0116】

ステップ2307においては、同一の物理記憶装置上に両データ構造が存在する部分に関してそれを分離するデータ移動案を作成する。ステップ2308においては、そのデータ移動案作成が成功したかどうか確認し、成功した場合にはステップ2310に進み、失敗した場合にはステップ2309に進みデータ再配置案作成処理失敗として処理を終了する。

【0117】

ステップ2310においては、作成されたデータ移動案を移動プラン情報750に記憶する。ステップ2311においては、作成されたデータ移動案に従ってデータ再配置ワーク情報670を更新し、ステップ2312に進む。

10

【0118】

ステップ2312においては、確認リストから現在確認対象となっているデータ構造の組のエントリを削除し、ステップ2303に進む。

【0119】

図20に、第3のデータ配置解析・データ再配置案作成処理である、データ構造の定義を基にした同時アクセス実行データ構造を分離するためのデータ再配置案作成処理の処理フローを示す。本処理においては、同時にアクセスされる可能性が高い、ログとその他のデータ、一時表領域とその他のデータ、表データとそれに対して作成された木構造の索引データが同一物理記憶装置18上に記憶されている部分が存在しないか確認をし、そのような部分が存在する場合にはそれを解決するデータ再配置案を作成する。

20

【0120】

ステップ2401で処理を開始する。ステップ2402では、DBMSデータ構造情報621を参照して全てのログであるデータ構造名561とそれを利用するDBMS110のDBMS名631の組を取得する。そして、そのデータ構造名とログ以外のデータを分離するデータ構造とする確認リストを作成し、ステップ2403に進む。

【0121】

ステップ2403ではステップ2402で作成した確認リストを用いてステップ2301から開始されるデータ構造分離のためのデータ再配置案作成処理を実行する。ステップ2404ではステップ2403におけるデータ再配置案作成処理が成功したか確認をし、成功した場合にはステップ2405に進む。失敗した場合にはステップ2412に進みデータ再配置案作成処理失敗として処理を終了する。

30

【0122】

ステップ2405では、DBMSデータ構造情報621を参照して全ての一時表領域であるデータ構造名561とそれを利用するDBMS110のDBMS名631の組を取得する。そして、そのデータ構造名と一時表領域以外のデータを分離するデータ構造とする確認リストを作成し、ステップ2406に進む。

【0123】

ステップ2406ではステップ2405で作成した確認リストを用いてステップ2301から開始されるデータ構造分離のためのデータ再配置案作成処理を実行する。ステップ2407ではステップ2406におけるデータ再配置案作成処理が成功したか確認をし、成功した場合にはステップ2408に進む。失敗した場合にはステップ2412に進みデータ再配置案作成処理失敗として処理を終了する。

40

【0124】

ステップ2408では、DBMS索引定義情報624を参照して全ての木構造索引の索引名635とそれに対応する表のデータ構造名を対応表情報637から取得する。そして、索引名635と対応する表のデータ構造名とそれらを保持するDBMS110のDBMS名631を組とする確認リストを作成し、ステップ2409に進む。

【0125】

ステップ2409ではステップ2408で作成した確認リストを用いてステップ2301

50

から開始されるデータ構造分離のためのデータ再配置案作成処理を実行する。ステップ 2410 ではステップ 2409 におけるデータ再配置案作成処理が成功したか確認をし、成功した場合にはステップ 2411 に進み、データ再配置案作成処理成功として処理を終了する。失敗した場合にはステップ 2412 に進みデータ再配置案作成処理失敗として処理を終了する。

【0126】

図 21 に第 4 のデータ配置解析・データ再配置案作成処理である、特定の表や索引の同一データ構造に対するアクセス並列度を考慮したデータ再配置案作成処理の処理フローを示す。この処理は、ランダムアクセス実行時の処理の並列度を考慮してディスクネックの軽減を図るためにデータの再配置を行うものである。この処理を実行する際には、データ再配置の確認対象とするデータ構造を DBMS 名 631 とデータ構造名 561 の組として指定する。

10

【0127】

ステップ 2501 で処理を開始する。ステップ 2502 において、指定されたデータ構造の物理記憶装置上に割り当てられた記憶領域利用総量を求める。この値は、指定データ構造がデータファイル上に割り当てられた容量と等しいため、DBMS データ記憶位置情報 622 中のデータ記憶位置を参照してその容量を求めればよい。

【0128】

ステップ 2503 においては、DBMS データ構造情報 621 を参照して指定データ構造における最大アクセス並列度 569 を取得する。ステップ 2504 において、ステップ 2502 で求めた指定データ構造の記憶領域利用総量をステップ 2503 で求めた最大アクセス並列度 569 で割った値を、指定データ構造の 1 つの物理記憶装置 18 上への割り当てを許可する最大量として求める。この制約により、特定の物理記憶装置 18 に偏ることなく最大アクセス並列度 569 以上の台数の物理記憶装置 18 に指定データ構造が分散して記憶されることになり、最大アクセス並列度 569 による並列度でランダムアクセスが実行されてもディスクネックになりにくい状況となる。なお、割り当て許可最大量の値は、実際のアクセス特性を考慮してこの方法で求めた値から更に増減させても構わない。

20

【0129】

ステップ 2505 において、指定データ構造のデータがステップ 2504 で求めた最大量を超えて 1 つの物理記憶装置 18 上に割り当てられているものが存在するかデータ再配置ワーク情報 670 を用いて確認し、そのようなものが存在しない場合にはステップ 2509 に進み、データ再配置案作成処理成功として処理を終了する。存在する場合にはステップ 2506 に進む。

30

【0130】

ステップ 2506 においては、ステップ 2504 で求めた最大量を超えて 1 つの物理記憶装置 18 上に割り当てられている部分を解消するデータ移動案を作成する。このとき、移動案作成に考慮するデータ移動量は指定データ構造の現在の物理記憶装置 18 上への割り当て量のステップ 2504 で求めた最大量からの超過分以上である必要がある。また、移動先物理記憶装置 18 においても、移動後にステップ 2504 で求めた最大量を超えないようにする必要がある。

40

【0131】

ステップ 2507 においては、ステップ 2506 のデータ移動案作成処理が成功したか確認をする。成功した場合にはステップ 2508 に進む。失敗した場合にはステップ 2510 に進み、データ再配置案作成処理失敗として処理を終了する。

【0132】

ステップ 2508 においては作成したデータ移動案を移動プラン情報 750 に記憶し、ステップ 2509 に進みデータ再配置案作成処理成功として処理を終了する。

【0133】

図 22 に第 5 のデータ配置解析・データ再配置案作成処理である、特定の表データに対するシーケンシャルアクセス時のディスクネックを解消するデータ再配置案作成処理の処理

50

フローを示す。この処理を実行する際には、データ再配置の確認対象とする表をDBMS名631とデータ構造名561の組として指定する。

【0134】

DBMS110毎により、シーケンシャルアクセス方法が定まっている。そこで、対象とするDBMS110の種類を絞り、あらかじめデータ位置管理主プログラム130がDBMS110におけるシーケンシャルアクセス方法を把握し、それらに対する最適化を行う。本実施の形態のDBMS110におけるシーケンシャルアクセス方法は以下の方法に従うものとする。あるデータ構造のデータをシーケンシャルアクセスする場合に、データ構造が記憶されているデータファイル名562とファイルブロック番号563を昇順にソートしその順序でアクセスを実行する。

10

【0135】

その他にシーケンシャルアクセス方法の決定方法としては、データファイルを管理する内部通番とファイルブロック番号563の組を昇順にソートした順番にアクセスする方法等が存在し、それらを利用したシーケンシャルアクセス方法の判断を実施してもよい。

【0136】

また、並列にシーケンシャルアクセスを実行する場合に、その領域の分割法は並列にアクセスしない場合のシーケンシャルにアクセスする順番を並列度に合わせて等分に分割するものとする。

【0137】

この並列アクセスによる分割後の1つのアクセス領域を全て同一の物理記憶装置18上に配置するのは必ずしも現実的ではない。そこで、分割後のアクセス領域がある一定量以上連続にまとまって1つの物理記憶装置上に記憶されていればよいと判断する。ただし、どのような場合でも連続してアクセスされることがなく、分割後のアクセス領域が異なるものに分類されるものに関しては、並列シーケンシャルアクセス時にアクセスがぶつかる可能性があるため、異なる物理記憶装置18に記憶するという指針を設けて、これに沿うようなデータ配置を作成することによりシーケンシャルアクセスの性能を高める。

20

【0138】

ステップ2601で処理を開始する。ステップ2602において、指定された表の物理記憶装置上に割り当てられた記憶領域利用総量を求める。この値は、指定データ構造がデータファイル上に割り当てられた容量と等しいため、DBMSデータ記憶位置情報622中のデータ記憶位置を参照してその容量を求めればよい。ステップ2603においては、DBMSデータ構造情報621を参照して指定データ構造における最大アクセス並列度569を取得する。

30

【0139】

ステップ2604において、ステップ2602で求めた指定表の記憶領域利用総量をステップ2603で求めた最大アクセス並列度569で割った量が、並列アクセス時にシーケンシャルにアクセスされる1つの領域のデータ量である。そこで、データ位置管理主プログラム130が把握している前述のシーケンシャルアクセス先の決定方法に基づいてDBMSデータ記憶位置情報622中の指定データ構造のデータファイルにおける記憶先を求め、それらのアクセス方法を前述のように予測し、その結果をもとに最大アクセス並列度569の並列アクセスが実行されると仮定した前述のデータ分割指針を作成する。

40

【0140】

ステップ2605において、データ再配置ワーク情報670を参照しながら、指定データ構造はステップ2604で作成したデータ分割指針に沿ったデータ配置が物理記憶装置18上で行われているか確認し、そうであればステップ2609に進み、データ再配置案作成処理成功として処理を終了する。そうでない場合にはステップ2606に進む。

【0141】

ステップ2606においては、物理記憶装置18上において、ステップ2604で求めたデータ分割指針に従うデータ配置を求める。このとき、データがある一定値以下の領域に細分化されている場合には、連続した空き領域を探し、そこにアクセス構造を保つように

50

データを移動するデータ移動案を作成する。また、最大アクセス並列度569の並列アクセスにより異なるアクセス領域に分離されるデータが同じ物理記憶装置18上に配置されないようなデータ移動案を作成する。

【0142】

ステップ2607においては、ステップ2606のデータ移動案作成処理が成功したか確認をする。成功した場合にはステップ2608に進み、失敗した場合にはステップ2610に進み、データ再配置案作成処理失敗として処理を終了する。

【0143】

ステップ2608においては、作成したデータ移動案を移動プラン情報750に記憶し、ステップ2609に進みデータ再配置案作成処理成功として処理を終了する。

10

【0144】

次に第6のデータ配置解析・データ再配置案作成処理である、特定のデータ構造に対する記憶装置10におけるキャッシュ効果を考慮したデータ再配置案作成処理の説明を行う。この処理を実行する際には、データ再配置の確認対象とするデータ構造としてDBMS名631とデータ構造名561を指定する。また、処理実行時に、記憶装置におけるキャッシュ効果が存在するかどうかを明示的に指定してもよい。前述のように、記憶装置構成情報134において、すべての記憶装置10がデータキャッシュ容量602に有効な値を保持しているわけではない。そのため、データキャッシュ容量602が無効値である記憶装置10は本処理の対象外とする。

【0145】

20

本処理においては、指定データ構造には記憶装置10のキャッシュ効果が存在するかどうかを判断する。まず、データ構造の単位データ量あたりのホストキャッシュにおける平均的キャッシュ利用量を計算する。その結果を用いて十分な量のホストキャッシュが利用可能かを判断する。十分な量のホストキャッシュを利用可能な場合には、アクセス頻度が低いデータに対してのみ記憶装置からデータを読み出すことになり、記憶装置におけるキャッシュ効果は極めて低いものとなる。

【0146】

その確認処理のフローを図23に示す。ステップ2801で処理を開始する。ステップ2802でDBMSキャッシュ構成情報625を参照して指定データ構造が属するキャッシュグループを求め、それからDBMSデータ記憶位置情報を参照してそのキャッシュグループに属するデータ構造の記憶のために割り当てられた領域量の総和を求める。

30

【0147】

ステップ2803において、DBMSキャッシュ構成情報625から指定データ構造が属するキャッシュグループに割り当てられたキャッシュサイズ566を求め、それとステップ2802で求めた領域量の総和から指定データ構造のホストにおける単位データ量あたりの平均キャッシュ利用量を求め、それをキャッシュ効果判断閾値と比較する。単位データ量あたりの平均キャッシュ利用量が閾値以上の場合にはステップ2804に進み、指定データ構造は記憶装置におけるキャッシュ効果がないと判断し、ステップ2806に進み処理を終了する。閾値未満の場合にはステップ2805に進み、指定データ構造は記憶装置におけるキャッシュ効果があると判断し、ステップ2806に進み処理を終了する。

40

【0148】

本処理において、記憶装置10におけるデータ構造のキャッシュ利用量を判断する。このとき、実際のデータのキャッシュ利用量はアクセスパターンに依存するが、ここでは平均的なケースを考え、データ構造のキャッシュ利用量は、その記憶装置10に割り当てられたデータ量に比例するとする。つまり、記憶装置構成情報134から記憶装置ボリューム物理記憶位置管理情報603を参照して記憶装置10における全データ記憶容量を求め、また、データキャッシュ容量602からデータキャッシュ容量を求める。これらの値から単位データ容量あたりのデータキャッシュ量が求まりこの値を元に判断する。また、記憶装置10において、提供するボリュームをいくつかグループ化し、それぞれで単位容量あたりのキャッシュ利用量を変化させる制御を行うことが可能である。この場合には、ボ

50

リユーム等の記憶領域毎にデータ構造のキャッシュ利用量が変化することになる。

【0149】

図24に特定のデータ構造に対する記憶装置10におけるキャッシュ効果を考慮したデータ再配置案作成処理の処理フローを示す。ステップ2701で処理を開始する。ステップ2702において、データ再配置案の作成対象として指定されたデータ構造に対して、処理開始時にキャッシュ効果があるかないかが明示的に指定されたか確認する。明示的に指定されていない場合にはステップ2703に進み、明示的に指定された場合にはステップ2704に進む。

【0150】

ステップ2703において、前述した指定データ構造の記憶装置10におけるキャッシュ効果が存在するかどうか確認処理を実行する。

10

【0151】

ステップ2704において、指定データ構造は記憶装置10においてキャッシュ効果があるかないかを確認する。キャッシュ効果があると判断された場合にはステップ2705に進み、ないと判断された場合にはステップ2706へ進む。

【0152】

ステップ2705においては、指定データ構造のデータを十分な量のキャッシュを利用可能な装置、ボリューム等の記憶領域に移動するデータ移動案を作成する。ここで、十分な量のキャッシュを利用可能な記憶領域とは、原則として単位容量あたりのキャッシュ利用量の大きな記憶領域を意味するが、以下の意味も併せ持つ。指定データ構造のホストキャッシュの利用可能量が大きなものである場合、アクセス頻度が高いものはホストキャッシュに留まることになる。そのため、記憶装置10におけるキャッシュ利用可能容量は、ホストキャッシュにおける利用可能容量に比べて、ある一定比率以上の容量がないとその効果は小さなものと考えられる。そこで、データ構造の単位データ量あたりのホストキャッシュにおける平均的キャッシュ利用量と記憶領域における単位容量あたりの平均的キャッシュ利用量の割合を計算し、その値と判断閾値を比較して記憶領域の方がよりキャッシュを多く利用できると判断される場合でないと十分な量のキャッシュが利用可能であると判断しない。本処理終了後、ステップ2707に進む。

20

【0153】

ステップ2706においては、指定データ構造のデータをキャッシュ利用量を少量に抑えられる領域にデータを移動する、つまり、記憶装置における単位容量あたりの平均的キャッシュ利用量が小さな装置、ボリューム等の領域へデータを移動するデータの移動案を作成し、ステップ2707に進む。

30

【0154】

ステップ2707においては、データ移動案の作成に成功したかどうかを確認する。成功した場合にはステップ2708に進み、失敗した場合にはステップ2710に進み、データ再配置案作成処理失敗として処理を終了する。

【0155】

ステップ2708においては作成したデータ移動案を移動プラン情報750に記憶し、ステップ2709に進みデータ再配置案作成処理成功として処理を終了する。

40

<第二の実施の形態>

本実施形態では、DBMSが実行される計算機と記憶装置が接続された計算機システムにおいて、データの記憶位置の管理を行う計算機が存在し、そこで計算機システム内のデータの記憶位置の管理を行う。計算機上で実行されるOS中のファイルシステムやボリュームマネージャは、動的にデータの記憶位置を変更する機能を有する。また、記憶装置においても、記憶装置内部でデータの記憶位置を動的に変更する機能を有する。

【0156】

データ記憶位置管理を実施する計算機は、DBMSに関する情報、データの記憶位置のマッピングに関する情報、記憶装置の構成情報を取得し、それらを用いて好適なデータ再配置案を作成する。ファイルシステム、ボリュームマネージャ、記憶装置に対して作成した

50

データ配置を指示し、それらのデータ再配置機能を用いてそのデータ再配置案に従ったデータ配置を実現する。

【0157】

図25は、本発明の第二の実施の形態における計算機システムの構成図である。図示されたように、本発明の第二の実施の形態は本発明の第一の実施の形態と以下の点が異なる。

【0158】

仮想ボリュームスイッチ72がI/Oパススイッチ72bに変更され、ネットワーク79とは未接続になる。DBホスト80bで実行されるOS100が有するボリュームマネージャ102がボリュームマネージャ102bに、ファイルシステム104がファイルシステム104bに変更され、OS100が保持するマッピング情報106がマッピング情報106bへとその内容が若干変更される。データ位置管理サーバ82内に記憶されるデータ記憶位置管理情報140がデータ記憶位置管理情報140内のホストマッピング情報650の内容を若干変更したホストマッピング情報650bに変更される。

10

【0159】

I/Oパススイッチ72bは仮想ボリュームスイッチ72と比較して経路制御の機能のみを保持する。本実施の形態においては、記憶装置10とDBホスト80b間のデータ転送を行うI/Oパス71とネットワーク79を異なるものとしているが、例えばiSCSIのような計算機と記憶装置間のデータ転送をネットワーク上で実施する技術も開発されており、本実施の形態においてもこの技術を利用してもよい。このとき、記憶装置10とDBホスト80bにおいてI/Oパスインターフェイス70が省かれ、計算機システム内からI/Oパス71とI/Oパススイッチ72bが省かれる構成となる。

20

【0160】

ボリュームマネージャ102bはボリュームマネージャ102と比べて、ボリュームマネージャ102bが提供する論理ボリュームの指定領域に記憶されているデータを、指定した記憶装置10が提供するボリュームの指定した記憶領域に移動する機能を有する。直接的にこの機能を有していない場合でも、管理領域の動的変更機能と管理領域内管理単位毎のデータ移動機能を有するものに関してはそれらの機能の組み合わせで実現できる。このデータ移動機能に関しては、ホスト上で実行される管理コマンドにより処理が実施される。

【0161】

ファイルシステム104bは、ファイルシステム104と比べて、ファイルシステム104bが管理する領域の中でデータが記憶されていない領域に対してファイルデータの一部をそこに移動する機能を有する。そのデータ移動指示方法は、移動するファイルとそのデータ領域と、移動先領域を指定するものとする。このデータ移動機能に関しては、ホスト上で実行される管理コマンドにより処理が実施される。

30

【0162】

ファイルのデータ記憶位置の動的変更機能は例えばファイルシステムのデフラグ機能として実現されており、その技術を用いてデータの移動先領域を指定できるようにすることにより前記データ移動機能を実現できる。データの移動先となることができる領域は、マッピング情報106bを参照することにより把握することができる。

40

【0163】

なお、本実施の形態においては、ボリュームマネージャ102bとファイルシステム104bのどちらか一方がデータ移動機能を有していれば良い。また、ファイルシステム104bがデータ移動機能を有している場合には、ボリュームマネージャ102bが存在しなくても本実施の形態に当てはめることができる。

【0164】

図26はOS100内に記憶されるマッピング情報106bを示す。図4のマッピング情報106からの変更点の概略は、マッピング情報106においては仮想ボリュームスイッチ72が提供する仮想ボリュームを利用していた部分をマッピング情報106bでは記憶装置10が提供するボリュームに変更された部分と、ファイル記憶位置情報530では保

50

持していなかったファイルシステムID535と空き領域の管理情報がファイル記憶位置情報530bに追加されたことである。

【0165】

マッピング情報106b中には、ボリュームローデバイス情報520b、ファイル記憶位置情報530bと論理ボリューム構成情報540bが含まれる。ボリュームローデバイス情報520b中にはOS100においてローデバイスを指定するための識別子であるローデバイスパス名521と、ローデバイスとして記憶装置10が提供するボリュームを利用する場合の記憶装置10の識別子である記憶装置名583、そのローデバイスによりアクセスされるボリュームあるいは論理ボリュームの識別子であるローデバイスボリューム名522bの組が含まれる。

10

【0166】

ファイル記憶位置情報540b中には、OS100においてファイルを指定するための識別子であるファイルパス名531とそのファイルが存在するファイルシステムのホスト内識別子であるファイルシステムID535とそのファイル中のデータ位置を指定するブロック番号であるファイルブロック番号532とそれに対応するデータが記憶されている記憶装置10が提供するボリュームもしくは論理ボリュームの識別子であるファイル配置ボリューム名533bと記憶装置10が提供するボリュームを利用する場合の記憶装置名583とそのボリューム上のデータ記憶位置であるファイル配置ボリュームブロック番号534の組が含まれる。ファイルパス名531が“Empty”であるエン트리536は特殊なエン트리であり、これはファイルシステム内における空き領域を示すデータである。この領域に対してデータの移動を行うことができる。

20

【0167】

論理ボリューム構成情報540b中にはボリュームマネージャ102bにより提供される論理ボリュームの識別子である論理ボリューム名541とその論理ボリューム上のデータの位置を示す論理ボリュームブロック番号542とその論理ブロックが記憶されている記憶装置10が提供するボリュームの識別子であるボリューム名501とそれを提供する記憶装置名583とボリューム上の記憶位置であるボリューム論理ブロック番号512の組が含まれる。

【0168】

図27はデータ位置管理サーバ82上に記憶されるホストマッピング情報650bを示す。図12のホストマッピング情報650からの変更点は、各ホストのマッピング情報106を保持するマッピング情報652が各ホストのマッピング情報106bを保持するマッピング情報652bになったことである。

30

【0169】

データ再配置処理において、データ移動機能を保持する部分が仮想ボリュームスイッチ72からボリュームマネージャ102b、ファイルシステム104bに変更されたことによる変更点は以下のようなものである。

【0170】

ボリュームマネージャ102b、ファイルシステム104bにおいては、DBホスト80b上において管理コマンドを実施することによりデータの移動を実施する。そこで、データ位置管理主プログラム130がネットワーク79を通してDBホスト80b上で実行されているデータ位置管理副プログラム120にボリュームマネージャ102b、ファイルシステム104bにおいてデータの移動を実施する管理コマンドを実行する指示を出し、それに従ってデータ位置管理副プログラム120が管理コマンドを実施することによりデータの移動を実施する。

40

【0171】

なお、ボリュームマネージャ102bにおいてデータ移動先として利用する領域は、データ位置管理主プログラム130移動指示を発行したときにはボリュームマネージャ102bの管理下に存在しない可能性がある。この場合には、データ位置管理副プログラム120は、データ移動の管理コマンドを実行する前に移動先領域をボリュームマネージャ10

50

2 bの管理下に収めるための管理コマンドを実行する。また、データ移動処理が終了後、他のDBホスト80上のボリュームマネージャ102bがデータ移動元の領域を新たなデータ移動先として利用するため、データ移動元領域の開放のための管理コマンドを実施する。

【0172】

また、データ再配置案作成処理においては、データ再配置案を作成する際に利用するワーク領域であるデータ再配置ワーク情報670がデータ再配置ワーク情報670bに、移動プラン情報750が移動プラン情報750bに変更される。

【0173】

図28にデータ再配置ワーク情報670bを示す。データ再配置ワーク情報670b中には、ワーク用記憶装置ボリューム記憶位置情報682とワーク用空き領域情報683とデータ構造仮想ボリューム内位置情報690bが含まれる。図14のデータ再配置ワーク情報670と比べて、仮想ボリューム物理記憶位置情報680はワーク用記憶装置ボリューム記憶位置情報682とワーク用空き領域情報683に分離・変更され、データ構造仮想ボリューム内位置情報690はデータ構造仮想ボリューム内位置情報690bに変更される。

10

【0174】

ワーク用記憶装置ボリューム記憶位置情報682は記憶装置10が提供するボリュームとその物理記憶装置18における記憶位置の一覧であり、記憶装置名583、ボリューム名501、ボリューム論理ブロック番号512、物理記憶装置名502、物理ブロック番号514の組を保持する。このデータは、記憶装置構成情報134を参照して初期化する。

20

【0175】

ワーク用空き領域情報683は、データ再配置案を作成する際にデータの移動先となりうる場所を管理するものであり、ホスト名631、ファイルシステムID535、論理ボリューム名541、論理ボリュームブロック番号542、記憶装置名583、ボリューム名501、ボリューム論理ブロック番号512、物理記憶装置名502、物理ブロック番号514の組のデータを保持する。このデータはワーク用空き領域情報683は、ホストマッピング情報650bと記憶装置構成情報134を参照して初期化する。このうち、ホスト名631とファイルシステムID535に有効値を有するエントリはそれらにより識別されるファイルシステム104bにおけるデータ移動先となりうる領域を示し、ホストマッピング情報650b中のファイル記憶位置情報530bからその領域を把握する。この領域に関しては、そのファイルシステム104b内に存在するデータの移動先として利用可能である。

30

【0176】

ホスト名631とファイルシステムID535に無効値を有するがボリューム名501に有効値を有するエントリは、まだどのホストからも利用されていない記憶装置10が提供しているボリューム内の記憶領域を示し、記憶装置構成情報134から把握される記憶装置10が提供している全領域からホストマッピング情報650bから把握される利用領域を除いた領域として把握する。この領域に関しては、ボリュームマネージャ102bにおけるデータ移動先として利用可能である。ボリューム名501に無効値を有するエントリは記憶装置10内のデータ移動先として利用可能な領域を示し、記憶装置構成情報134から把握可能である。

40

【0177】

データ構造仮想ボリューム内位置情報690bはDBMS110が保持しているデータ構造が記憶装置10が提供するボリュームのどこに記憶されているかを示した情報であり、ホスト名651、DBMS名631、データ構造名561、データファイルパス名562、ファイルシステムID535、ファイルブロック番号563、記憶装置名583、ボリューム名501、ボリューム論理ブロック番号512の組を保持する。この情報は、DBMSスキーマ情報136中のDBMSデータ記憶位置情報622とDBMSホスト情報626とホストマッピング情報650bを参照し、ファイル(ローデバイス)パス、(論理

50

) ボリュームに関して対応する部分をまとめることにより初期データを作成する。

【0178】

図29に移動プラン情報750bを示す。移動プラン情報750bは、ボリュームマネージャ102bに対するデータ移動指示を記憶する論理ボリューム移動プラン情報753とファイルシステム104bに対するデータ移動指示を記憶するファイルブロック移動プラン情報754と記憶装置10に対するデータ移動指示を記憶する物理記憶位置移動プラン情報752を含む。これら情報に関しては、何もデータを持たないように初期化する。図15の移動プラン情報750と比較して、移動プラン情報750bは仮想ボリューム移動プラン情報752が削除され、論理ボリューム移動プラン情報753とファイルブロック移動プラン情報754が追加されたものとなっている。

10

【0179】

論理ボリューム移動プラン情報753には、移動順序761、移動処理を行うホスト名631、移動元ボリュームとその領域を指定する移動論理ボリューム名773、移動論理ボリュームブロック番号774、移動先の記憶装置10とその記憶領域を指定する移動先記憶装置名764、移動先ボリューム名765、移動先ブロック番号766の組が記憶される。ファイルブロック移動プラン情報754には、移動順序761、移動処理を行うホスト名651、移動元ファイルとその領域を指定するファイルシステムID535、移動データファイルパス名775、移動ファイルブロック番号776、移動先領域を指定する移動先記憶装置名764、移動先ボリューム名765、移動先ブロック番号766の組が記憶される。

20

< 第三の実施の形態 >

本実施形態では、DBMSが実行される計算機とファイルを管理単位とする記憶装置がネットワークを用いて接続された計算機システムにおいて、データの記憶位置の管理を行う計算機が存在し、そこで計算機システム内のデータの記憶位置の管理を行う。計算機上で実行されるOS中のネットワークファイルシステムは、複数のファイルを仮想的な1つのファイルにまとめ、その構成を動的に変更する機能を有する。また、記憶装置においても、記憶装置内部でデータの記憶位置を動的に変更する機能を有する。

【0180】

データ記憶位置管理を実施する計算機は、DBMSに関する情報、データの記憶位置のマッピングに関する情報、記憶装置の構成情報を取得し、それらを用いて好適なデータ再配置案を作成する。ネットワークファイルシステム、記憶装置に対して作成したデータ配置を指示し、それらのデータ再配置機能を用いてそのデータ再配置案に従ったデータ配置を実現する。

30

【0181】

図30は、本発明の第三の実施の形態における計算機システムの構成図である。図示されたように、本発明の第三の実施の形態は本発明の第一の実施の形態と以下の点が異なる。

【0182】

本実施の形態においてはI/Oパスインターフェイス70、I/Oバス71、仮想ボリュームスイッチ72が存在せず、記憶制御装置10cとDBホスト80cはネットワーク79を介してのみ接続される。記憶装置10はファイルを単位としたデータ記憶管理を行う記憶装置10cに変更される。そのため、物理記憶装置稼動情報32、ボリューム物理記憶位置管理情報36がそれぞれ物理記憶装置稼動情報32c、ファイル記憶管理情報36cに変更される。

40

【0183】

DBホスト80cで実行されるOS100ではボリュームマネージャ102、ファイルシステム104が削除されその代わりにネットワークファイルシステム104cが追加され、OS100が保持するマッピング情報106がマッピング情報106cへ変更される。データ位置管理サーバ82内に記憶される記憶装置稼動情報132、記憶装置構成情報134、データ記憶位置管理情報140がそれぞれ記憶装置稼動情報132c、記憶装置構成情報134c、データ記憶位置管理情報140内のホストマッピング情報650の内容

50

を若干変更したホストマッピング情報 6 5 0 c に変更される。

【 0 1 8 4 】

記憶装置 1 0 はファイルを管理単位とする記憶装置 1 0 c に変更される。D B ホスト 8 0 c からのアクセスも N F S 等のファイルをベースとしたプロトコルで実施される。記憶装置 1 0 におけるボリュームの役割は、記憶装置 1 0 c においてはファイルもしくはファイルを管理するファイルシステムとなり、そのファイルの記憶位置管理情報がファイル記憶管理情報 3 6 c である。1 つの記憶装置 1 0 c の中に複数のファイルシステムが存在しても構わない。

【 0 1 8 5 】

物理記憶装置 1 8 の稼動情報はボリュームを単位とした取得からファイルシステム又はファイル単位とした取得に変更する。記憶装置 1 0 c 内にファイルシステムが存在する場合でもデータの移動機能を実現可能であり、データ移動指示方法は、前と同じく移動するファイルとそのデータ領域と、移動先領域を指定するものとする。本実施の形態においては、記憶装置 1 0 c におけるデータ移動機能は必須であるとする。

10

【 0 1 8 6 】

ネットワークファイルシステム 1 0 4 c は、記憶装置 1 0 c が提供するファイルをアクセスするための機能を提供する。更に、複数のファイルを仮想的な 1 つのファイルとして提供する。この機能を実現するためにネットワークファイルシステム 1 0 4 c は管理情報をマッピング情報 1 0 6 c 内に保持し、仮想ファイルアクセス時にこの管理情報を参照し、実際のアクセス先を求める処理を行う。更に、その構成を動的に変更する機能を有する。これら処理は、D B ホスト 8 0 上で管理コマンドを実行することにより実施される。

20

【 0 1 8 7 】

図 3 1 は記憶装置 1 0 c 内に保持される物理記憶装置稼動情報 3 2 c を示す。図 2 の物理記憶装置稼動情報 3 2 からの変更点は、稼動情報取得単位がボリュームからファイルシステムに変更されたため、ボリューム名 5 0 1 の部分がファイルシステム名 1 0 0 1 に変更されたことである。また、稼動情報取得単位をファイルとしてもよく、このときはボリューム名 5 0 1 の部分がファイルシステム名 1 0 0 1 とファイルパス名 1 0 0 2 に変更される。

図 3 2 は記憶装置 1 0 c 内に保持されるファイル記憶管理情報 3 6 c を示す。図 3 のボリューム物理記憶位置管理情報 3 6 からの変更点は、ボリューム物理記憶位置メイン情報 5 1 0、ボリュームデータ移動管理情報 5 1 1 からファイル物理記憶位置情報 5 1 0 c、ファイルデータ移動管理情報 5 1 1 c にそれぞれ変更される。上記の変更内容は、ボリュームの識別子がボリューム名 5 0 1 がファイルの識別子となるファイルシステム名 1 0 0 1 とファイルパス名 1 0 0 2 に、ボリューム内のデータ領域を示すボリューム論理ブロック番号 5 1 2 と移動論理ブロック番号 7 8 2 がそれぞれファイルブロック番号 1 0 0 3 または移動ファイルブロック番号 1 0 2 1 に変更されたことである。

30

【 0 1 8 8 】

ここで、ファイルパス名 1 0 0 2 が “ E m p t y ” であるエントリ 1 0 1 5 は特殊なエントリであり、このエントリには記憶装置 1 0 c 内の物理記憶装置 1 8 の領域のうち、指定ファイルシステム内でファイルの記憶領域としてに割り当てられていない領域を示し、図 3 中のボリュームデータ移動管理情報 5 1 1 を用いるデータ移動方式で説明した処理手順を用い、この領域に対して移動するデータをコピーすることによりデータの物理記憶位置の動的変更機能を実現する。

40

【 0 1 8 9 】

図 3 3 は D B ホスト 8 0 c の O S 1 0 0 内に記憶されているマッピング情報 1 0 6 c を示す。マッピング情報 1 0 6 c 中には、ネットワークファイルシステムマウント情報 1 0 3 0 と仮想ファイル情報 1 0 4 0 と仮想ファイルデータ移動管理情報 1 0 5 0 が含まれる。

【 0 1 9 0 】

ネットワークファイルシステムマウント情報 1 0 3 0 は、記憶装置 1 0 c から提供され、D B ホスト 8 0 c においてマウントされているファイルシステムの情報で、ファイルシス

50

テムの提供元記憶装置とそのファイルシステムの識別子である記憶装置名583とファイルシステム名1001、そして、そのファイルシステムのマウントポイントの情報であるマウントポイント1031の組を保持する。

【0191】

仮想ファイル情報1040は、ネットワークファイルシステム104cが提供する複数の記憶装置から提供されるファイルを仮想的な1つのファイルとして提供する機能の管理に用いる情報で、提供される仮想ファイルの識別子である仮想ファイルパス名1041とそのデータ領域を示す仮想ファイルブロック番号1042とそのデータ領域のデータを実際に保持するファイルの識別子である構成ファイルパス名1043とその記憶領域を示す構成ファイルブロック番号1044の組を含む。

10

【0192】

仮想ファイルデータ移動管理情報は1050は、ネットワークファイルシステム104cが提供する仮想ファイルの構成の変更処理の一部である、構成データの移動処理を行う際に利用する管理情報で、データ移動を行う移動元の仮想ファイルの識別子である移動仮想ファイルパス名1051とデータ移動を行う領域を示す移動仮想ファイルブロック番号1052とそのデータの移動先ファイルとその移動先領域を示す移動先構成ファイルパス名1053と移動先ファイルブロック番号1054、そして、データ移動処理を行う際の管理情報である差分管理情報785とコピーポイント786の組を含む。データの移動先に関しては、移動先ファイルの移動先指定領域に記憶領域実体が確保されている以外の制約は存在しない。本情報を利用し、図3中のボリュームデータ移動管理情報511を用いるデータ移動方式で説明した処理手順を用いることにより、データの移動機能を実現できる。

20

【0193】

図34はデータ位置管理サーバ82c上に記憶される記憶装置稼働情報132cを示す。図8の記憶装置稼働情報132からの変更点は、記憶装置稼働情報132ではボリューム単位で稼働情報を取得していたものを記憶装置稼働情報132cではファイルシステム単位で稼働情報を取得していることである。そのため、ボリューム名501がファイルシステム名1001に変更される。

【0194】

図35はデータ位置管理サーバ82c上に記憶される記憶装置構成情報134cを示す。図9の記憶装置構成情報134からの変更点は、記憶装置毎のボリューム物理記憶位置メイン情報510を記憶する記憶装置ボリューム物理記憶位置管理情報604が記憶装置毎のファイル物理記憶位置情報510cを記憶する記憶装置ファイル物理記憶位置情報604cに変更される。

30

【0195】

図36はデータ位置管理サーバ82c上に記憶されるホストマッピング情報650cを示す。図12のホストマッピング情報650からの変更点は、各ホストのマッピング情報106を保持するマッピング情報652が各ホストのマッピング情報106cを保持するマッピング情報652cになったことである。

【0196】

データ再配置処理において、データ移動機能を保持する部分が仮想ボリュームスイッチ72からネットワークファイルシステム104cに変更され、記憶装置10cがファイルを管理単位とするように変更されたことによる変更点は以下のようなものである。

40

【0197】

ネットワークファイルシステム104cにおいては、DBホスト80c上において管理コマンドを実施することによりデータの移動を実施する。そこで、データ位置管理主プログラム130がネットワーク79を通して、DBホスト80c上で実行されているデータ位置管理副プログラム120に、ネットワークファイルシステム104cにおいてデータの移動を実施する管理コマンドを実行する指示を出し、それに従ってデータ位置管理副プログラム120が管理コマンドを実施することによりデータの移動を実施する。このとき、

50

現在空きの領域に対するデータの移動を行うため、データ移動先領域となるファイルやファイル内の領域が存在しない。このようなデータ移動指示を受け取った場合には、ネットワークファイルシステム 104c は指定ファイルの新規ファイル作成や領域拡張を実施し、それが成功した後にデータ移動処理を開始し、データ移動中の領域不足の問題を回避する。

【0198】

また、ネットワークファイルシステム 104c は通常のプロトコルを利用してファイル作成や領域拡張を実施する。そのため、そのデータ記憶先が必ずしも最適な場所に割り当てられるとは限らない。そこで、ネットワークファイルシステム 104c によるデータ移動が完了後、今度は記憶装置 10c に対して記憶装置内データ移動の指示を出し、作成したデータ再配置案に従ったデータ配置を実現する。このとき、記憶装置 10c 内でデータ移動元と移動先が重なっている場合には、一旦データを移動先とは異なる空き領域にデータを移動し、その後再度指定された移動先にデータを移動させる処理を記憶装置 10c は実行する。

10

【0199】

データ再配置案作成処理においては、データ再配置案を作成する際に利用するワーク領域であるデータ再配置ワーク情報 670 がデータ再配置ワーク情報 670c に、移動プラン情報 750 が移動プラン情報 750c に変更される。

【0200】

図 37 にデータ再配置ワーク情報 670c を示す。データ再配置ワーク情報 670c 中には、記憶装置ファイル物理記憶位置情報 681c とデータ構造記憶装置内ファイル位置情報 690c が含まれる。図 14 のデータ再配置ワーク情報 670 と比べて、仮想ボリューム物理記憶位置情報 680 は記憶装置ファイル物理記憶位置情報 681c に変更され、データ構造仮想ボリューム内位置情報 690 はデータ構造仮想ボリューム内位置情報 690c に変更される。

20

【0201】

記憶装置ファイル物理記憶位置情報 681c は記憶装置 10c が提供するファイルシステムとその内部に存在するファイル、その物理記憶装置 18 における記憶位置の一覧であり、記憶装置名 583、ファイルシステム名 1001、ファイルパス名 1002、ファイルブロック番号 1003、物理記憶装置名 502、物理ブロック番号 514 の組を保持する。このデータは、記憶装置構成情報 134c を参照して初期化する。ファイルパス名 1002 が "Empty" であるエントリ 1071 は記憶装置 10c のファイルシステム名 1001 の領域のうち、ファイルの記憶に利用されていない領域を示し、ここに対してデータの移動が可能である。

30

【0202】

データ構造仮想ボリューム内位置情報 690c は DBMS 110 が保持しているデータ構造が記憶装置 10c が提供するファイルのどこに記憶されているかを示した情報であり、ホスト名 651、DBMS 名 631、データ構造名 561、データファイルパス名 562、ファイルブロック番号 563、記憶装置名 583、ファイルシステム名 1001、ファイルパス名 1002、ファイルブロック番号 1003 の組を保持する。この情報は、DBMS スキーマ情報 136 中の DBMS データ記憶位置情報 622 と DBMS ホスト情報 626 とホストマッピング情報 650c を参照し、ファイルパスに関して対応する部分をまとめることにより初期データを作成する。

40

【0203】

図 38 に移動プラン情報 750c を示す。移動プラン情報 750c は、ネットワークファイルシステム 104c に対するデータ移動指示を記憶する仮想ファイルブロック移動プラン情報 755 と記憶装置 10c に対するデータ移動指示を記憶する物理記憶位置移動プラン情報 752c を含む。これら情報に関しては、何もデータを持たないように初期化する。図 15 の移動プラン情報 750 と比較して、移動プラン情報 750c は仮想ボリューム移動プラン情報 752 が削除され、仮想ファイル移動プラン情報 755 が追加され、物理

50

記憶位置移動プラン情報 7 5 2 が物理記憶位置移動プラン情報 7 5 2 c に変更されたものとなる。

【 0 2 0 4 】

仮想ファイル移動プラン情報 7 5 5 には、移動順序 7 6 1、移動処理を行うホスト名 6 5 1、移動元の仮想ファイルとその領域を指定する移動仮想ファイルパス名 1 0 5 1、移動仮想ファイルブロック番号 1 0 5 2、移動先の構成ファイルとその領域を指定する移動先構成ファイルパス名 1 0 5 3、移動先ファイルブロック番号 1 0 5 4 の組が記憶される。物理記憶位置移動プラン情報 7 5 2 c には、移動順序 7 6 1、移動処理を行う移動記憶装置名 7 6 7、移動元のファイルとその領域を指定する移動ファイルシステム名 1 1 0 1、移動ファイルパス名 1 1 0 2、移動ファイルブロック番号 1 1 0 3、移動先の物理記憶装置 1 8 とその領域を指定する移動先物理記憶装置名 7 7 1、移動先物理ブロック番号 7 7 2 の組が記憶される。

10

【 0 2 0 5 】

本実施例においては、ネットワークファイルシステム 1 0 4 c が複数のファイルから 1 つの仮想ファイルを構成する機能を有しているが、ここで、単純に 1 つの仮想ファイルは 1 つのファイルから構成されても構わない。このとき、ネットワークファイルシステム 1 0 4 c は動的なデータ移動機能のみを提供する。さらに、DBMS 1 1 0 が処理を中断しても構わない場合には、DBMS 1 1 0 の処理の中断後、ファイルのコピーを実施し、アクセス先がコピー先になるようにシンボリックリンクを張りなおした後に DBMS 1 1 0 が処理を再開することによるデータ再配置も可能である。さらに、ネットワークファイルシステム 1 0 4 c におけるデータ移動を行わず、記憶装置 1 0 c のみでデータ移動を実施することも可能である。

20

【 0 2 0 6 】

【発明の効果】

本発明により以下のことが可能となる。第一に、DBMS が管理するデータの特性を考慮して記憶装置におけるデータ記憶位置を決定することにより、DBMS に対してより好ましいアクセス性能特性を持つ記憶装置を保持する計算機システムが実現される。これにより、その計算機システムで稼動している DBMS の性能を向上させることができる。特に、複数の記憶装置を利用する DB システムにおいて、各記憶装置間へアクセス要求が適切に分散化されるようにし、DBMS の処理性能を向上させる。

30

【 0 2 0 7 】

第二に、DBMS が稼動している計算機システムにおいて、DBMS の特性を考慮した記憶装置の良好なアクセス性能特性達成を目的としたデータ記憶位置再配置処理を実現するため、計算機システムの性能に関する管理コストを削減することができる。特に、本発明を用いることにより、データ記憶位置の再配置案を自動で作成することが可能であり、管理コストの削減に大きく寄与する。更に、多数の DBMS が稼動し、多数の記憶装置が存在するシステムにおいても利用可能で集中的な管理を可能とするため、そのようなシステムにおける性能に関する管理コストを削減することができる。

【図面の簡単な説明】

【図 1】第一の実施の形態における計算機システムの構成を示す図である。

40

【図 2】記憶装置 1 0 内に保持されている物理記憶装置稼動情報 3 2 を示す図である。

【図 3】記憶装置 1 0 内に保持されているボリューム物理記憶位置管理情報 3 6 を示す図である。

【図 4】第一の実施の形態における DB ホスト 8 0 の OS 1 0 0 内に記憶されているマッピング情報 1 0 6 を示す図である。

【図 5】DBMS 1 1 0 内に記憶されているその内部で定義・管理しているデータその他の管理情報であるスキーマ情報 1 1 4 を示す図である。

【図 6】DB ホスト 8 0 のメモリ 8 8 上に記憶されている実行履歴情報 1 2 2 を示す図である。

【図 7】仮想ボリュームスイッチ 7 2 が保持する仮想ボリューム情報 7 3 を示す図である

50

- 。
- 【図 8】データ位置管理サーバ 8 2 上に記憶される記憶装置稼動情報 1 3 2 を示す図である。
- 【図 9】データ位置管理サーバ 8 2 上に記憶される記憶装置構成情報 1 3 4 を示す図である。
- 【図 1 0】データ位置管理サーバ 8 2 上に記憶される D B M S スキーマ情報 1 3 6 を示す図である。
- 【図 1 1】データ位置管理サーバ 8 2 上に記憶される D B M S 実行履歴情報 1 3 8 を示す図である。
- 【図 1 2】データ位置管理サーバ 8 2 上に記憶されるデータ記憶位置管理情報 1 4 0 を示す図である。 10
- 【図 1 3】データ位置管理主プログラム 1 3 0 によるデータ再配置処理の処理フローを示す図である。
- 【図 1 4】データ配置解析・再配置案作成処理で利用するデータ再配置ワーク情報 6 7 0 を示す図である。
- 【図 1 5】データ配置解析・再配置案作成処理により作成されるデータ移動案を格納する移動プラン情報 7 5 0 を示す図である。
- 【図 1 6】記憶装置稼動情報 1 3 2 を基にした同時アクセス実行データ構造を分離するためのデータ再配置案作成処理の処理フローを示す図である。
- 【図 1 7】D B M S 実行履歴情報 1 3 8 を基にした同時アクセス実行データ構造を分離するためのデータ再配置案作成処理で利用するクエリ実行時同時アクセスデータカウント情報 7 0 0 を示す図である。 20
- 【図 1 8】D B M S 実行履歴情報 1 3 8 を基にした同時アクセス実行データ構造を分離するためのデータ再配置案作成処理処理フローを示す図である。
- 【図 1 9】指定されたデータ構造とそのデータ構造と同時にアクセスされる可能性が高いデータ構造の組を分離するデータ再配置案を作成する処理のフローを示す図である。
- 【図 2 0】データ構造の定義を基にした同時アクセス実行データ構造を分離するためのデータ再配置案作成処理の処理フローを示す図である。
- 【図 2 1】特定の表や索引の同一データ構造に対するアクセス並列度を考慮したデータ再配置案作成処理の処理フローを示す図である。 30
- 【図 2 2】特定の表データに対するシーケンシャルアクセス時のディスクネックを解消するデータ再配置案作成処理の処理フローを示す図である。
- 【図 2 3】特定のデータ構造に対する記憶装置 1 0 におけるキャッシュ効果を考慮したデータ再配置案作成処理で利用するキャッシュ効果判定処理の処理フローを示す図である。
- 【図 2 4】特定のデータ構造に対する記憶装置 1 0 におけるキャッシュ効果を考慮したデータ再配置案作成処理の処理フローを示す図である。
- 【図 2 5】第二の実施の形態における計算機システムの構成を示す図である。
- 【図 2 6】D B ホスト 8 0 b の O S 1 0 0 内に記憶されるマッピング情報 1 0 6 b を示す。
- 。
- 【図 2 7】データ位置管理サーバ 8 2 上に記憶されるホストマッピング情報 6 5 0 b を示す図である。 40
- 【図 2 8】データ配置解析・再配置案作成処理で利用するデータ再配置ワーク情報 6 7 0 b を示す図である。
- 【図 2 9】データ配置解析・再配置案作成処理により作成されるデータ移動案を格納する移動プラン情報 7 5 0 b を示す図である。
- 【図 3 0】第三の実施の形態における計算機システムの構成を示す図である。
- 【図 3 1】記憶装置 1 0 c 内に保持される物理記憶装置稼動情報 3 2 c を示す図である。
- 【図 3 2】記憶装置 1 0 c 内に保持されるファイル記憶管理情報 3 6 c を示す図である。
- 【図 3 3】D B ホスト 8 0 c の O S 1 0 0 内に記憶されているマッピング情報 1 0 6 c を示す図である。 50

【図34】データ位置管理サーバ82c上に記憶される記憶装置稼動情報132cを示す図である。

【図35】データ位置管理サーバ82c上に記憶される記憶装置構成情報134cを示す図である。

【図36】データ位置管理サーバ82c上に記憶されるホストマッピング情報650cを示す図である。

【図37】データ配置解析・再配置案作成処理で利用するデータ再配置ワーク情報670cを示す図である。

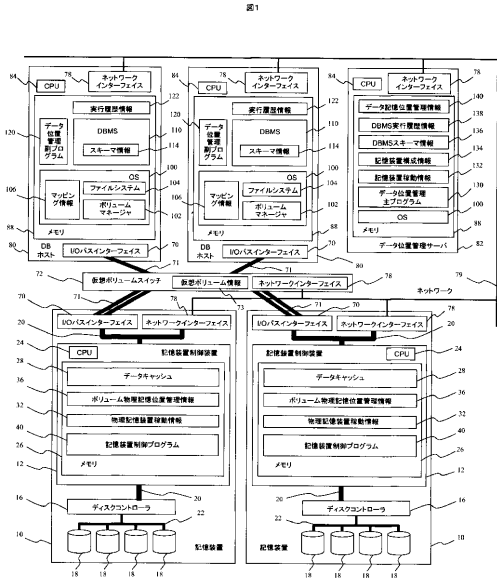
【図38】データ配置解析・再配置案作成処理で利用する移動プラン情報750cを示す図である。

10

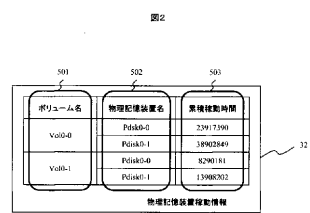
【符号の説明】

10, 10c	記憶装置	
18	物理記憶装置	
28	データキャッシュ	
32, 32c	物理記憶装置稼動情報	
36	ボリューム物理記憶位置管理情報	
36c	ファイル記憶管理情報	
70	I/Oパスインターフェイス	
71	I/Oパス	
72	仮想ボリュームスイッチ	20
72b	I/Oパススイッチ	
73	仮想ボリューム情報	
78	ネットワークインターフェイス	
79	ネットワーク	
80, 80b, 80c	DBホスト	
82	データ位置管理サーバ	
100	OS (オペレーティングシステム)	
102, 102b	ボリュームマネージャ	
104, 104b	ファイルシステム	
104c	ネットワークファイルシステム	30
106, 106b, 106c	マッピング情報	
110	DBMS (データベース管理システム)	
114	スキーマ情報	
120	データ位置管理副プログラム	
122	実行履歴情報	
130	データ位置管理主プログラム	
132, 132c	記憶装置稼動情報	
134, 134c	記憶装置構成情報	
136	DBMSスキーマ情報	
138	DBMS実行履歴情報	40
140	データ記憶位置管理情報	
650, 650b, 650c	ホストマッピング情報	

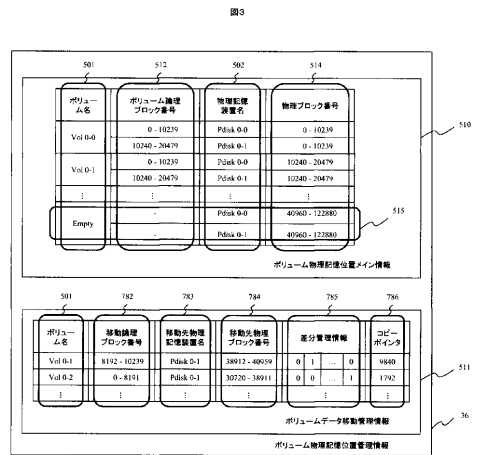
【図1】



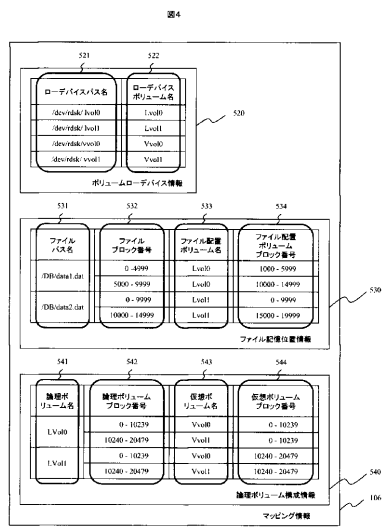
【図2】



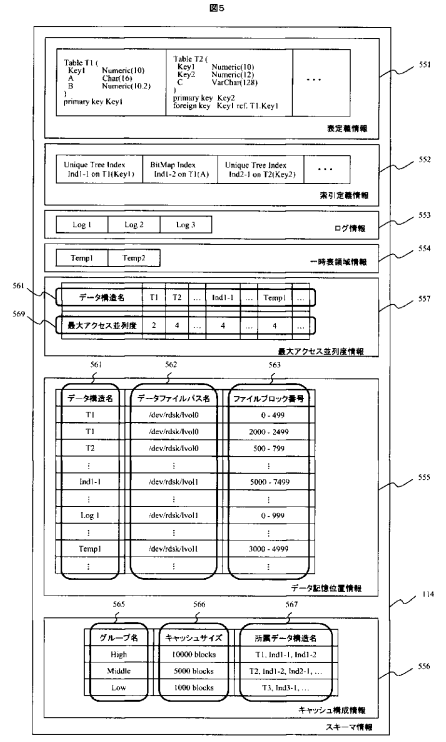
【図3】



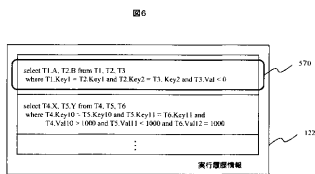
【図4】



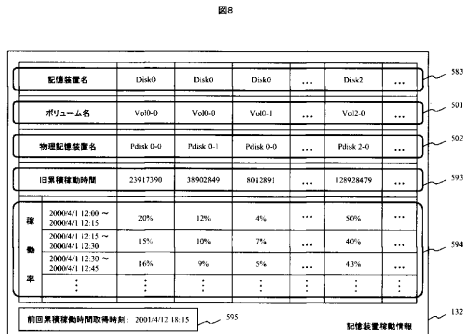
【図5】



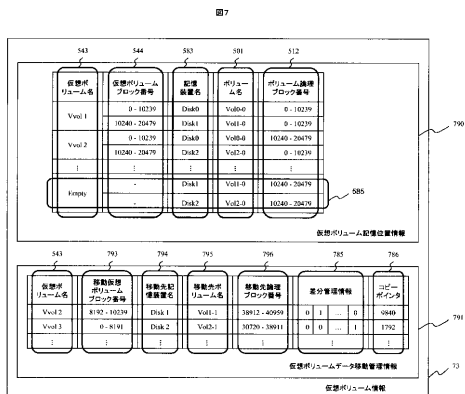
【図6】



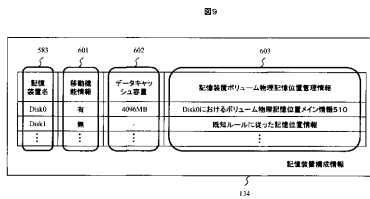
【図8】



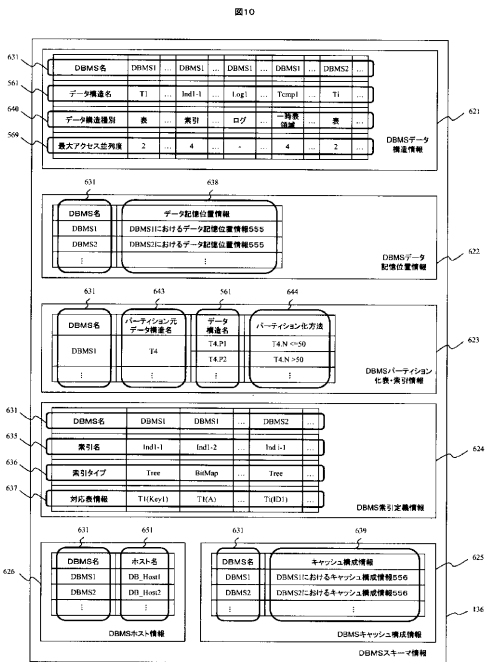
【図7】



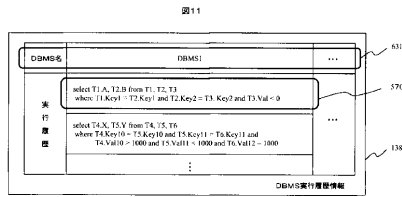
【図9】



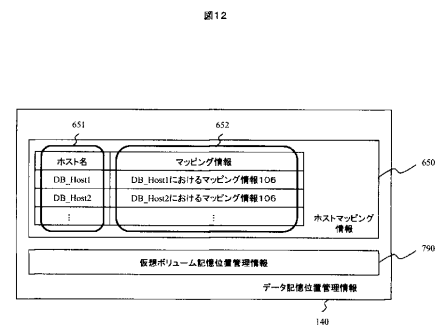
【図10】



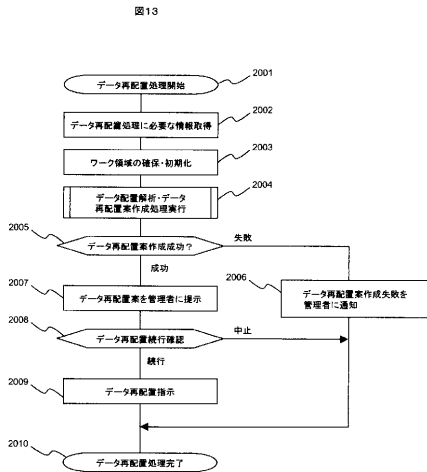
【図11】



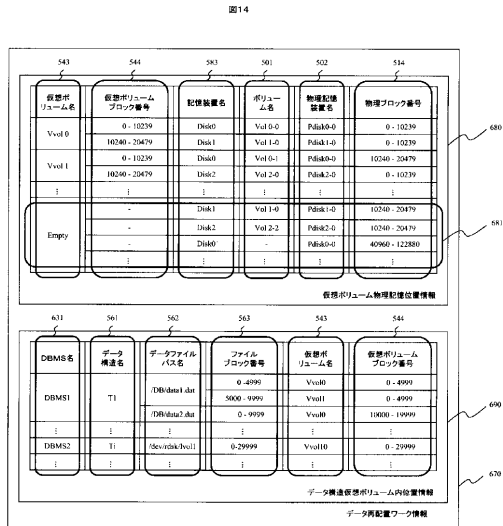
【図12】



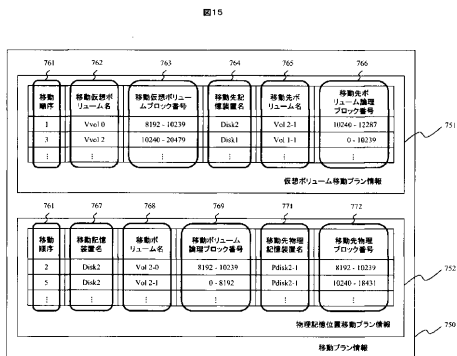
【図13】



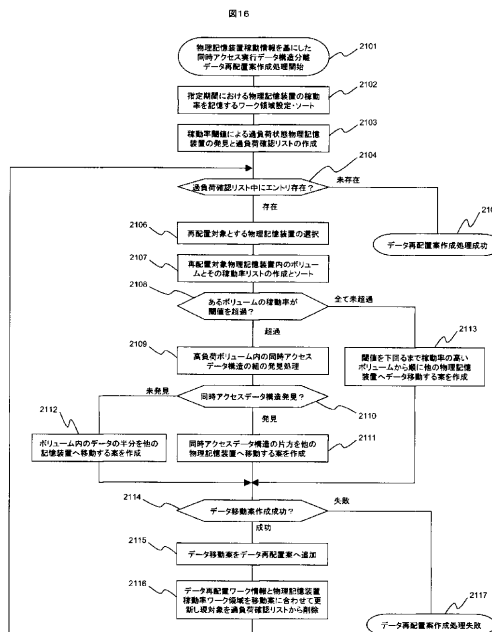
【図14】



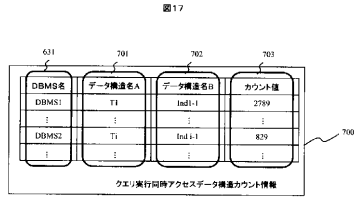
【図15】



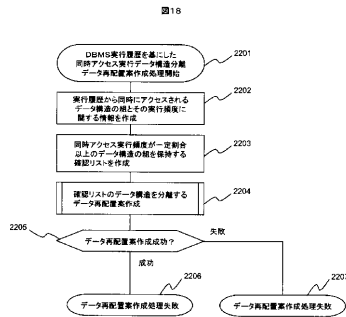
【図16】



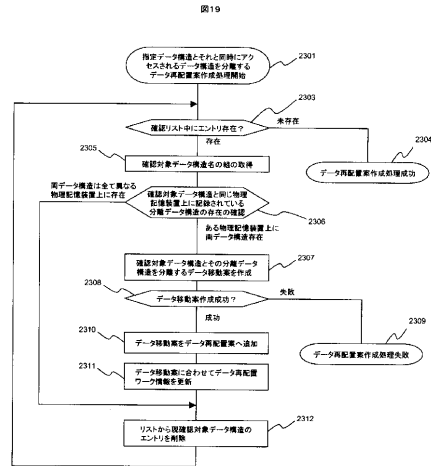
【図17】



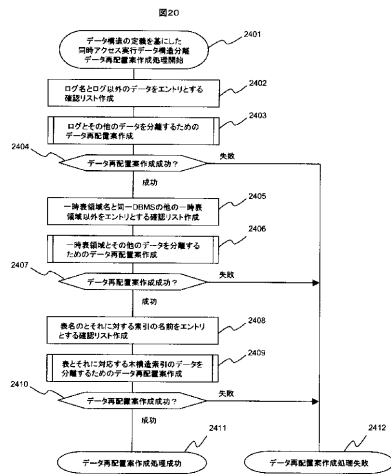
【図18】



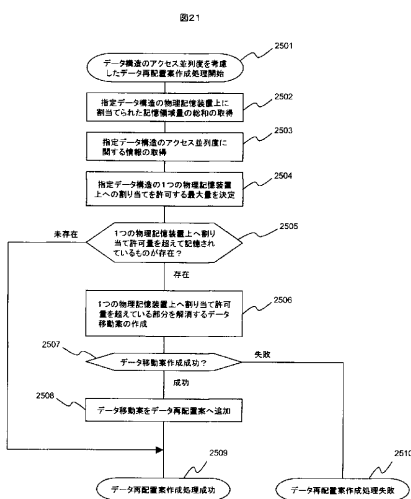
【図19】



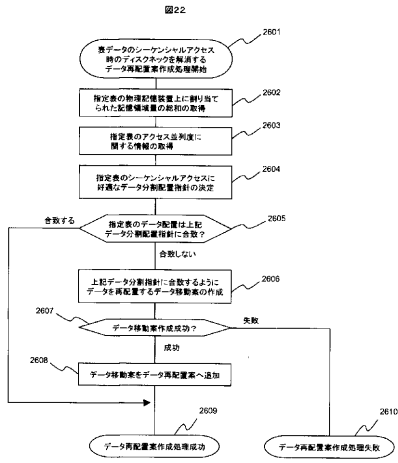
【図20】



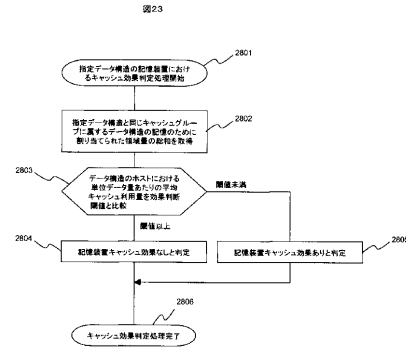
【図21】



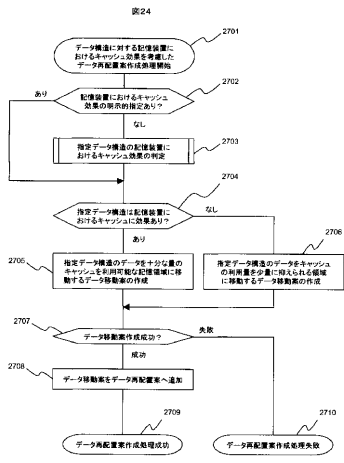
【図22】



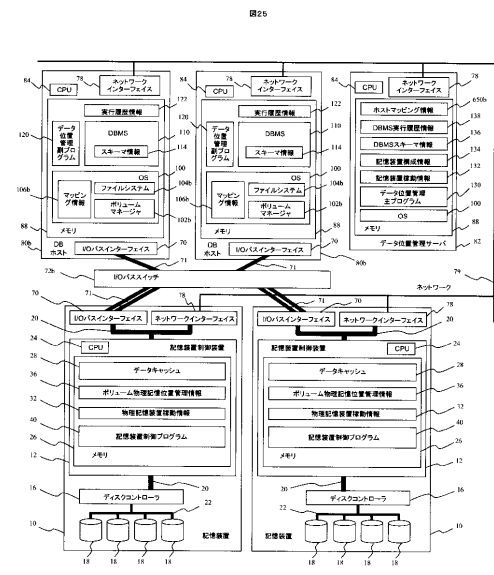
【図23】



【図24】

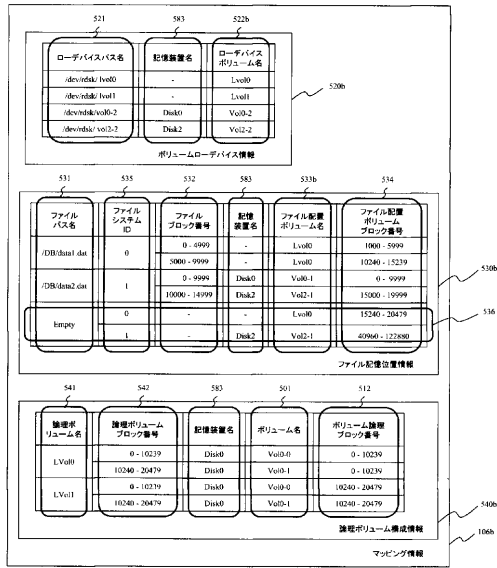


【図25】



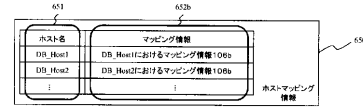
【図26】

図26



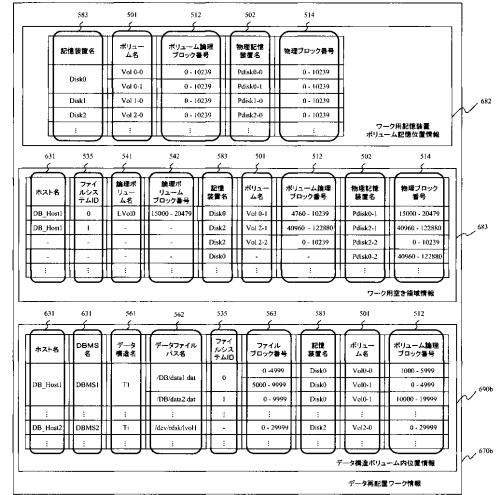
【図27】

図27



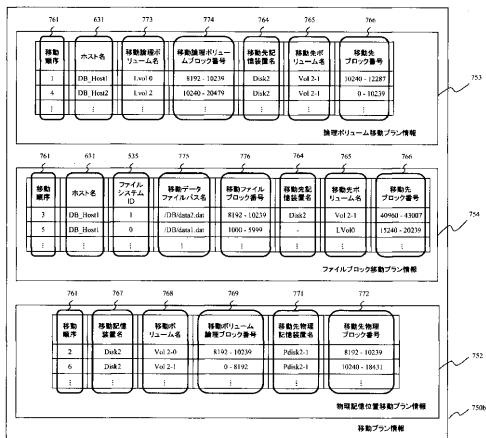
【図28】

図28



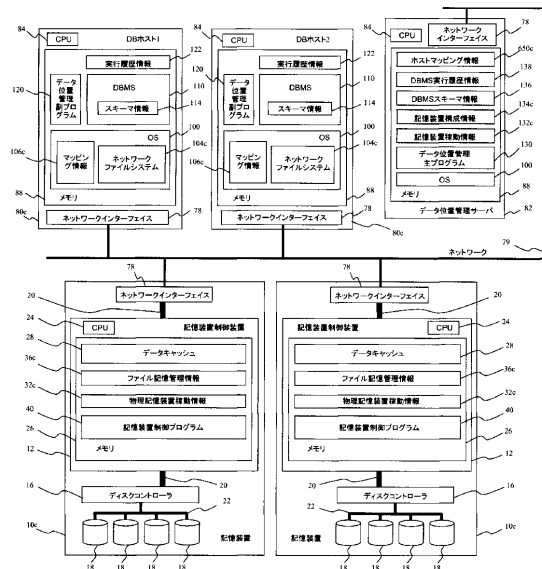
【図29】

図29

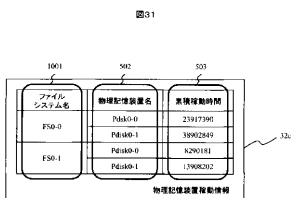


【図30】

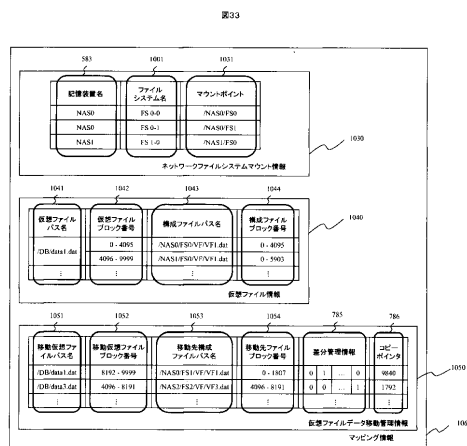
図30



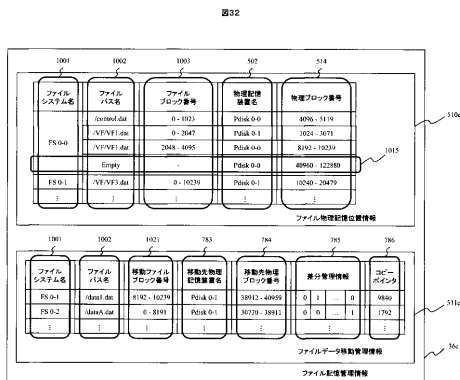
【図31】



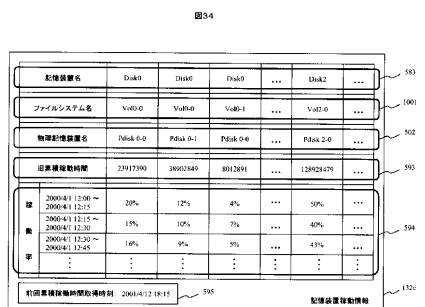
【図33】



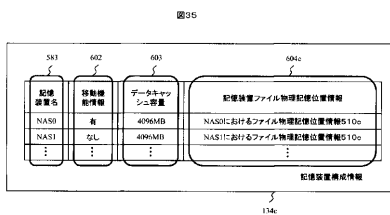
【図32】



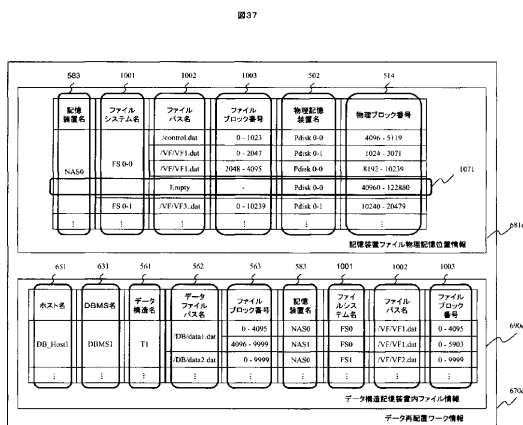
【図34】



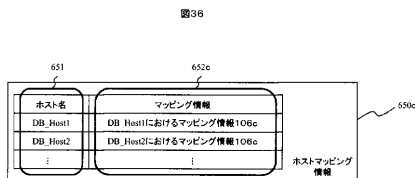
【図35】



【図37】

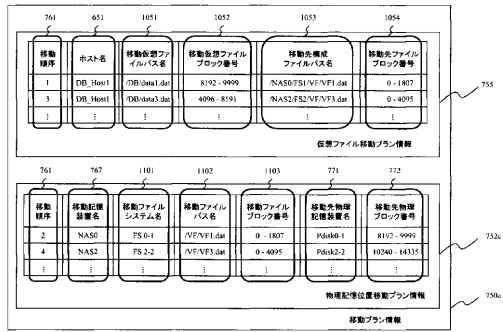


【図36】



【図38】

図38



フロントページの続き

(72)発明者 喜連川 優
千葉県松戸市二十世紀が丘丸山町17

審査官 高瀬 勤

(56)参考文献 特開平05-012338(JP,A)
原 憲宏 他, 並列DBサーバシステムにおけるDB再配置方式の提案, 第48回(平成6年前期)全国大会講演論文集(4), 日本, 社団法人情報処理学会, 1994年3月25日, p.207-208

(58)調査した分野(Int.Cl., DB名)

G06F 12/00

G06F 17/30

JSTPlus(JDream2)